



King's Research Portal

DOI:

[10.1017/S1755020316000435](https://doi.org/10.1017/S1755020316000435)

Document Version

Peer reviewed version

[Link to publication record in King's Research Portal](#)

Citation for published version (APA):

Nicolai, C. (2017). Equivalences for truth predicates. *Review Of Symbolic Logic*, 10(2), 322-356.
<https://doi.org/10.1017/S1755020316000435>

Citing this paper

Please note that where the full-text provided on King's Research Portal is the Author Accepted Manuscript or Post-Print version this may differ from the final Published version. If citing, it is advised that you check and use the publisher's definitive version for pagination, volume/issue, and date of publication details. And where the final published version is provided on the Research Portal, if citing you are again advised to check the publisher's website for any subsequent corrections.

General rights

Copyright and moral rights for the publications made accessible in the Research Portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognize and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the Research Portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the Research Portal

Take down policy

If you believe that this document breaches copyright please contact librarypure@kcl.ac.uk providing details, and we will remove access to the work immediately and investigate your claim.

THE REVIEW OF SYMBOLIC LOGIC
Volume 0, Number 0, Month 2009

EQUIVALENCES FOR TRUTH PREDICATES

CARLO NICOLAI

Munich Center for Mathematical Philosophy
LMU Munich

Abstract. One way to study and understand the notion of truth is to examine principles that we are willing to associate with truth, often because they conform to a pre-theoretical or to a semi-formal characterization of this concept. In comparing different collections of such principles, one requires formally precise notions of inter-theoretic reduction that are also adequate to compare these conceptual aspects. In this work I study possible ways to make precise the relation of conceptual equivalence between notions of truth associated with collections of principles of truth. In doing so, I will consider refinements and strengthenings of the notion of relative truth-definability proposed by (Fujimoto 2010): in particular I employ suitable variants of notions of equivalence of theories considered in (Visser 2006; Friedman & Visser 2014) to show that there are better candidates than mutual truth-definability for the role of sufficient condition for conceptual equivalence between the semantic notions associated with the theories. In the concluding part of the paper, I extend the techniques introduced in the first and show that there is a precise sense in which ramified truth (either disquotational or compositional) does not correspond to iterations of comprehension.

§1. The Project Comparing theories is one of the fundamental activities of the scientist and of the philosopher. Often this comparison is carried out via formalization, and there is a great deal of controversy on how to properly formalize scientific or mathematical theories.¹ In this work I focus on a narrower target and investigate notions of reduction and equivalence involving deductive systems obtained by extending a base arithmetical theory with principles of truth. This description is deliberately vague at this stage: I will in fact relate either two theories of truth or a theory of truth and an extension of the base theory with suitable comprehension schemata. In the rest of this introduction, I will motivate the project in its twofold nature.

1.1. Truth and Conceptual Equivalence Recent works on truth and semantic paradoxes fall in two broad categories. On one side we find works that aim at describing a class of models for languages endowed with a truth predicate. These construction are carried out by defining the notion of truth in a more powerful metatheory – usually a system of analysis in the case of an arithmetical base theory. We call this approach the *semantic* approach (for an overview, see (Field 2008)). On the other side authors that give priority to rules and principles characterizing a primitive notion of truth. On this approach one focuses on formal systems that generate truths (and falsities). The latter approach is widely known as the *axiomatic* approach ((Halbach 2014; Horsten 2011)).

¹ See (Lutz 2016) for an up to date account of the debate in the context of philosophy of science.

That there can be fruitful interactions between the two *modus operandi* is clear since Tarski's seminal work.² Collections of principles and rules usually form succinct and accessible bases for motivating, evaluating, criticizing some approaches to semantic paradoxes. Semantic constructions are often indispensable in the creative act of articulating one such approach.³

In order to characterize the notion of truth arising from a semantic construction, one usually fixes a specific model – e.g. the minimal fixed point of Kripke's theory in Strong Kleene logic from (Kripke 1975) – and reads off the properties of the truth predicate in this model. If one, by contrast, gives priority to the analysis of systems generating truths, the vicinity to a semantic construction may be part of the characterization of the system's truth predicate, but also other features of this truth predicate need to be part of the analysis. For instance, the vicinity of a system of truth T to a class of models Σ may be evaluated, as proposed by (Fischer et alii 2015), by considering ω -models. On this approach, T will capture Σ if and only if⁴

$$(1) \quad (\mathbb{N}, S) \models T \Leftrightarrow (\mathbb{N}, S) \text{ belongs to } \Sigma$$

where S is the extension of the truth predicate.

However, given some S, Σ , there are several examples of systems T_1, T_2 satisfying (1) that cannot be plausibly associated to the same notion of truth.⁵ This suggests that, in order to characterize the notion of truth associated with a system of truth, we need to integrate the vicinity to a semantic construction – that may still not necessarily be identified with the criterion of sameness of ω -models of (Fischer et alii 2015) – with further, complementary criteria. A list of such criteria for systems of truth will be provided in §2.

Before sketching the project any further, a note on terminology. We will employ interchangeably the expressions 'notion of truth', 'concept of truth', and 'conception of truth' in relation to the outcome of the reflection on the properties of a truth predicate as depicted

² Cf. (Fischer et alii 2015) and (Halbach 2014) for an in-depth treatment of the relationships between the two approaches.

³ A Kantian aphorism by Hannes Leitgeb may help in depicting the situation: it is often the case that axiomatic truth without semantic truth is empty, whereas semantic truth without axiomatic truth is blind.

⁴ For simplicity, we consider theories in classical logic, so \models is the classical satisfaction relation.

⁵ For instance, let us consider the theories extending a base syntax theory B – that can be taken to be a weak arithmetical system such as EA considered below – with typed Tarski biconditionals of the form

$$' \varphi ' \text{ is true if and only if } \varphi$$

with φ a sentence of \mathcal{L}_B , and the theory obtained by extending B with full, compositional axioms of the form

$$\text{for all sentences } \varphi, \psi \text{ of } \mathcal{L}_B, \text{ the sentence ' } \varphi \text{ and } \psi \text{ ' is true if and only if ' } \varphi \text{ ' is true and ' } \psi \text{ ' is true.}$$

The two theories are different in many respects: the latter theory is capable of proving many more general claims than the first, it is compositional and not disquotational as the former, and its notion of truth is not reducible to the logical resources of B . Moreover, the compositional theory will *not* in general be relatively interpretable in B , for finitely axiomatizable B , and not conservative over it if the right amount of inductive reasoning with truth is available. Nonetheless, the two truth theories have the same models based on \mathbb{N} – that is a pair (\mathbb{N}, S) where S is the unique set of truths for \mathcal{L} – so they capture the same semantic construction according to the criterion above and, if the 'reductionist' picture is endorsed, their notions of truth coincide.

by a theory of truth. In all cases we *intend a cluster of conditions satisfied by the truth predicates of the languages of our theories*. In this sense a notion (conception, concept) of truth is primarily not a psychological notion, and so my focus diverges from the classical ontological and scientific questions on concepts that may be found in the literature.⁶

At the informal level, we want to address the question: when the concepts of truth can be considered equivalent? Still intuitively, our guiding answer is: when they equally possess or lack a number of crucial properties. This characterization is of course still unsatisfactory. In the first place we need the additional assumption that conceptions of truth are given by formal theories, and formal systems in particular. Moreover, the relevant properties that two theories of truth may or may not possess are dependent on one's stance on truth: an instrumentalist may regard provability of truth-free consequences as the only criterion for equivalence of two truth predicates, whereas someone who is interested in the purely truth theoretic principles of the theories may completely disregard their consistency strength. In addition, one would like to integrate the mutual preservation of some features of theories of truth qua mathematical objects into a more general framework of equivalence of formal systems. It would be at least puzzling to realize that systems T_1 and T_2 are formally equivalent in a strong sense but their truth predicates fare differently with respect to some of the characterizing, quasi-formal desiderata generally imposed on notions of truth.

In the first part of this work I will therefore investigate the possibility of finding a formal counterpart to the informal relation of 'conceptual equivalence' just sketched. This is not an easy nor an original task: some authors, notably (Fujimoto 2010) and (Halbach 2014), have tried to come up with notions of reductions up to the task. In particular, by applying the notion of *relative truth-definability* – that is a relative interpretation that keeps the syntactic vocabulary unchanged – and variations thereof, they aimed at capturing the notion of *conceptual equivalence* of truth predicates (see §1 of (Fujimoto 2010)), or comparing the *conceptual strength* of two (or more) truth predicates. In the first part of this paper I will show that the notions of (proof-theoretic reductions) that are commonly employed in the literature, including relative truth-definability, do not suffice to adequately characterize the relation of 'conceptual equivalence of truth predicates'. In the second part I will consider stricter notions of reduction that may represent a sufficient condition for two theories' truth predicates to be conceptually equivalent, although arguably not a necessary one due to their strictness.

1.2. Classes and Truths In formal terms, the result of extending a base arithmetical theory with truth axioms has been often regarded as another, perhaps more succinct way of extending it with further ontological assumptions on the existence of sets of natural numbers. Famously (Feferman 1991) has employed truth axioms as a device to investigate the limits of predicativity given the natural numbers. This programme finds its roots in the inter-reducibility of suitable truth axioms with certain fragments of ramified analysis. These mutual reductions may be of interest for multiple reasons. In general, relating truth with certain forms of membership may help in harmonizing analyses of semantic, set-theoretic, and property-theoretic paradoxes. Moreover, trading-off ontological commitment to sets with semantic commitment to truth axioms may be attractive for philosophers

⁶ For various approaches to the notion of concept in philosophy and cognitive science, we refer to (Laurence & Margolis 1999). (Woodfield 1991) also considers conceptions, but his notion is very close to the usual notion of concept employed in psychology.

interested in reducing assumptions on the existence of sets to syntactic objects and ideology (Halbach 2014).⁷

The strict notions of reduction that will be investigated in the first part of the paper will shed light on the folklore reductions between suitable truth axioms and certain forms of comprehension – the ones considered by Halbach and Feferman, for instance – and yield a surprising picture in which semantic and set-theoretic assumptions *fail to be equivalent*. In particular, in this picture ontological commitments seem to be reducible to semantic assumptions but not vice versa.

§2. Truth predicates and theoretical equivalence In this section I extract from the recent literature some desiderata for adequacy of systems of truth (I mainly refer to (Leitgeb 2007; Halbach & Horsten 2015)) to find plausible candidates for conceptual properties of a theory's truth predicate besides the vicinity to some semantic construction. The criteria seem to oscillate between genuinely truth theoretic, or 'conceptual' ones, and criteria on the systems *qua* mathematical objects, although in the axiomatic context the boundaries between the two kinds are often very difficult to trace. In the following list, *T* and *W* will denote theories extending a base theory *B* with one or more truth predicates. The fundamental idea behind the listing is that it cannot be the case *T* and *W* are seen as conceptually equivalent and yet they disagree on some of the desiderata below. It should be clear the list is by no means exhaustive: I simply appeal to criteria that acquired some general consensus to guide our formal analysis.

Ontological commitments. This requirement can be paraphrased as follows: suppose *T* and *W* display the same conception of truth. If *T* enables one to interpret the objects of the domain to which truth is ascribed as the intended bearers of truth, typically sentence types, also *W* should do so. A more general formulation of this criterion may be as follows: sameness of conception entails preservation of the ontological commitments of the theories.⁸ The very possibility of comparing truth predicates, let alone equating them, appears to be rooted in the possibility of applying them to structurally similar linguistic-mathematical objects. It is known, for instance, that there are theories of truth that do not have ω -models, such as FS (cf. (Halbach 2014; Friedman & Sheard 1987)). It would be embarrassing, for any satisfactory notion of equivalence of conception of truth, to pair FS with theories that admit standard models.

Truth-theoretic generalizations. The provability of generalizations involving the truth predicate may be seen as a formal rendering of the requirement, for the truth predicate, of

⁷ For instance, in (Halbach 2014), we read:

...I expect that not only truth theories can be applied to eliminate ontological commitments, but that the work on proof theory also sheds light on how semantic and ontological commitment are related. [...] In particular, the reductions of second-order theories to truth theories may be taken as evidence that ontological commitment can be replaced with ideological or semantic commitment, or perhaps even that there is no very clear distinction between the two kinds of commitment. (p. 318)

⁸ Preservation of ontological commitment may also be intended not as *identity* of the syntactic universe, but as isomorphism of the structures of the 'objects of truth'. I will allow for this possibility, although I will not focus on similar cases.

adequately ‘expressing infinite conjunctions’ or allowing for ‘semantic ascent’.⁹ In other words if the notions of truth associated with T and W are conceptually equivalent, we expect the existence of an effective procedure that enables one to translate (provable) generalizations involving the truth predicate(s) of T , e.g. stating that all members of a countable set \mathcal{S} of formulas are true, into general claims involving *the truth predicate(s)* of W , in which the formulas expressing membership in \mathcal{S} in W and T apply to the same syntactic objects. The requirement in *italic* should not be underestimated: we will consider below cases of theories that satisfy the condition just stated only by means of resources other than truth.

Reducibility. The debate over truth-theoretic deflationism has prompted an extensive literature on the reducibility of semantic resources to the underlying object theory – for an overview, cf. (Horsten 2011). There are at least three notions of reductions that may be considered. The first is definability: it is obvious that this is not available in the case of truth. The second is conservativity: a notion of truth may help in reaching consequences that were not provable by the object theory only. The third is relative interpretability: although the notion of truth is not definable in the syntactic base theory, its behaviour may be replicated by notions that work on a suitable relativization of the syntactic domain. Conservativity and interpretability, besides having deep philosophical implications, are measures of strength of the theories of truth. We require them to be preserved in adequate formal renderings of conceptual equivalence of truth predicates.

Compositionality. In general, we welcome the possibility of understanding the truth of a compound sentence only by knowing the truth values of the compounding sentences (for as many objects to which the truth predicate can be applied as we can find). For the truth predicates of T and W to display the same conception of truth, both theories should have this property: if $Tr_T(\varphi \wedge \psi)$ can be understood – in T – as ‘ $Tr_T \ulcorner \varphi \urcorner$ and $Tr_T \ulcorner \psi \urcorner$ ’, also $Tr_W(\ulcorner \varphi \wedge \psi \urcorner)^*$ (where $(\cdot)^*$ is an appropriate isomorphism between the syntactic universes), should be understood analogously in W . It cannot be the case, for instance, that T is able to decompose sentences compositionally only for an initial portion of the syntactic universe (perhaps only for standard syntactic objects), whereas W can do so for *any* sentence in the domain of our quantifiers. As a consequence, this criterion enables us to separate between disquotational and compositional theories.

Symmetry. Under the assumption that T features one single truth predicate Tr , the internal theory of T is defined as the set of φ such that $T \vdash Tr \ulcorner \varphi \urcorner$. The external theory of T simply the set of φ such that $T \vdash \varphi$. When the external and the internal theories of T coincide, T is said to be symmetric. If T and W embody equivalent conceptions of truth and if T can be consistently made symmetric, this should also be possible for W . A mismatch between the internal and the external theories, although often unavoidable, is an unpleasant sign: it may suggest that the meaning externally assigned to logical connectives shifts once we move under the scope of the truth predicate.

Finite axiomatizability (with no additional resources). Often accompanied by the compositionality requirement, the finite axiomatizability of a theory of truth is sometimes related to the learnability, by human beings, of truth values for infinitely many sentences starting with finite instructions.¹⁰ It would appear at least questionable for T and W being

⁹ Some authors think that equating ‘expressing infinite conjunctions’ with ‘proving universally quantified statements involving truth’ is a hasty move. This is the case of (Halbach 1999), for instance. Some others, such as (Cieśliński 2010), find the equation harmless.

¹⁰ This is a famous Davidsonian theme. See (Davidson 1984).

considered conceptually equivalent, but T be presented as a finite set of clauses whereas W as an infinite list of instructions. This desideratum is often stressed in relation to Davidson's program, but it is not clear why the requirement cannot be extended to *recursive*, instead of finite, sets of formulas. After all, isn't the schema $\phi \rightarrow \phi$ learnable, even though it has infinitely many instantiations? Perhaps a distinction can be made by instances of schemata characterizing logical vocabulary only and instances of schemata concerning nonlogical vocabulary, but it's not our intention here to follow this line of reasoning any further. Since we are dealing with formal systems of truth, in fact, there is enough evidence supporting the claim that the truth predicate of a finitely axiomatizable theory may considerably differ from the truth predicate of a non finitely axiomatizable theory; for instance, the former may not be reducible to the object theory, whereas the latter may well be.¹¹

Type restrictions. We will require the (non) applicability of the truth predicate to sentences containing semantic vocabulary to be part of the conception of truth arising from some axiomatic theory. That said, it is not immediately clear what the distinction between typed and type-free truth may be.¹² We shall be content with the following condition: T will be dubbed as type-free if it proves the truth of a *particular* sentence containing the truth predicate.

2.1. Theories, Base Theories, Theories of Truth Theories of truth require a theory of the objects to which truth applies: in our case eternal, context-free sentence-types. In turn it is well-known that, modulo synonymy (cf. §2.2.), direct axiomatizations of syntax or finite sets coincide with suitable arithmetical theories. The main role of the underlying syntax theory is to state the relevant syntactic notions and operations needed for truth ascriptions and to prove their relevant properties: to this purpose, the standard choice of PA as base theory seems unnecessarily strong. In this work I will employ elementary arithmetic EA as base theory: it is a properly weaker theory than PA or $\text{I}\Sigma_1$.¹³ EA captures in fact in a direct and natural way the standard development of the syntax of first-order theories, as carried out for instance in (Feferman 1960; Smorynski 1977) for PA and PRA.¹⁴ More on the choice of EA can be found in §4.1..

Officially, EA is formulated in a first-order, relational language \mathcal{L} with logical constants in $\{\neg, \wedge, \forall\}$ featuring finitely many relation symbols $Z(x)$ (' x is identical to zero'), $S(x, y)$ (' y is the successor of x '), $E(x, y)$ (' y is 2^x '), $A(x, y, z)$ (' z is the sum of x and y '), $M(x, y, z)$ (' z is the product of x and y '). With $x \neq y$, expressions of the form $(\forall x \leq y)(\phi(x))$ and $(\exists x \leq y)(\phi(x))$ are said to be obtained from $\phi(v)$ by *bounded quantification*. The class of formulas of \mathcal{L} that contain only bounded quantifiers will be referred to as the class of *elementary* formulas. The set of axioms of EA contains, besides the logical axioms of pred-

¹¹ We discuss the finite axiomatizability criterion further in §4.1..

¹² Ch. 10 proposes the satisfaction of the criteria as sufficient conditions for a system of truth T to be typed:

1. Only sentences not containing used or mentioned occurrences of the truth predicate should be deemed true by the theory;
2. the theory should not rule out the possibility of picking any set definable in the base theory as extension of the truth predicate restricted to sentences containing the truth predicate itself.

¹³ More precisely, the consistency of EA can be proved already in PRA, which is known to be conservatively extended by $\text{I}\Sigma_1$ for Π_2 -formulas.

¹⁴ Yet, we are far from being close to the theoretical lower bounds. For this purpose, one might choose a theory interpretable in Q as in (Nicolai 2016).

icate logic with identity formulated in a Hilbert-style calculus, the functionality (including totality) axioms for the relation symbols just introduced and axioms corresponding to their recursive clauses such as (the universal closures of)

$$\begin{aligned} Z(y) \wedge S(y, z) &\rightarrow (E(x, z) \leftrightarrow x = y) \\ S(x_0, x_1) \wedge A(z_0, z_0, z_1) &\rightarrow (E(x_0, z_0) \leftrightarrow E(x_1, z_1)) \end{aligned}$$

In addition, we have the schema of bounded induction for elementary (Δ_0) formulas:

$$(\text{Ind-}\Delta_0) \quad Z(x) \wedge \varphi(x) \wedge \forall y, y_0 (\varphi(y) \wedge S(y, y_0) \rightarrow \varphi(y_0)) \rightarrow \forall y \varphi(y)$$

Partial truth-definitions are available in IA_0 : they can be used to show that EA is finitely axiomatizable (cf. Ch. V). In working with EA, I will often employ functional expressions for elementary functions,¹⁵ such as $S(x)$, as unofficial counterparts of the relations just introduced. In practice I work in a definitional extension of EA, assuming in the background the possibility of translating back the unofficial abbreviations into the official signatures by eliminating terms via (suitably bounded) existential quantifiers.

The formalization of the syntax of first-order theories in EA is carried out without difficulties: the details of one such coding can be found, for instance, in (Schwichtenberg and Wainer 2012). Par abus de language, we denote with $\ulcorner e \urcorner$ the formal code of e . We also avail ourselves with symbols for elementary functions such that EA ‘proves’ the following equations, with φ, ψ \mathcal{L}_0 -formulas and x_i a (meta)variable of \mathcal{L}_0 – standing for v_{i1} :

$$\begin{aligned} R(\ulcorner x_1 \urcorner, \dots, \ulcorner x_n \urcorner) &= \ulcorner R(x_1, \dots, x_n) \urcorner & \text{ng}(\ulcorner \varphi \urcorner) &= \ulcorner \neg \varphi \urcorner & \text{dn}(\ulcorner \varphi \urcorner) &= \ulcorner \neg \neg \varphi \urcorner \\ \text{and}(\ulcorner \varphi \urcorner, \ulcorner \psi \urcorner) &= \ulcorner \varphi \wedge \psi \urcorner & \text{all}(\ulcorner v \urcorner, \ulcorner \varphi \urcorner) &= \ulcorner \forall v \varphi \urcorner \end{aligned}$$

The operation $\varphi(v), t \mapsto \varphi(t)$ of substituting a ‘term’ for a free variable in a formula is naturally represented in EA by an elementary formula $\text{sub}(x, y, u, v)$ such that

$$\begin{aligned} \text{EA} \vdash \forall x, y, u \exists! v \text{sub}(x, y, u, v) \\ \text{EA} \vdash \text{sub}(\ulcorner \varphi(x_i) \urcorner, \ulcorner x_i \urcorner, \ulcorner t \urcorner, y) \rightarrow y = \ulcorner \varphi(t) \urcorner \end{aligned}$$

For notational convenience, I will write sub as if it were a function. We also write $\langle x_1, \dots, x_n \rangle$ to designate in EA a finite sequence, $(x)_i$ to indicate the i^{th} element of the sequence x , $\text{lh}(x)$ for its length. All these operations correspond to elementary functions and are provably total in EA. I will also extensively employ the following elementary syntactic notions

¹⁵ The class of *elementary functions* \mathcal{E} is obtained by closing the initial functions $\text{zero}(\cdot)$, $\text{suc}(\cdot)$, $+$, \times , 2^x , $P_i^n(x_1, \dots, x_n) = x_i$ with $(1 \leq i \leq n)$, truncated subtraction $x \dot{-} y$ under the operations of composition and bounded minimalization:

$$H(\vec{x}) = F(G_1(\vec{x}), \dots, G_n(\vec{x})); \quad (\mu t \leq y) P(\vec{x}, t) = \begin{cases} \text{the least } t \leq y \text{ s.t. } P(\vec{x}, t) \\ 0, & \text{if there is no such } t \end{cases}$$

where F, G are elementary functions and P an elementary predicate. EA has sufficient resources to naturally introduce new relations corresponding to the elementary functions by proving their defining equations.

relative to some elementary presented theory T in \mathcal{L}_0 :

| | |
|-----------------------------------|---|
| $\text{Fml}_{\mathcal{L}_0}(x)$ | ‘ x is the code of an \mathcal{L}_0 -formula’ |
| $\text{Fml}_{\mathcal{L}_0}^i(x)$ | ‘ x is the code of an \mathcal{L}_0 -formula with i variables free’ |
| $\text{Sent}_{\mathcal{L}_0}(x)$ | ‘ x is the code of an \mathcal{L}_0 -sentence’ |
| $\text{Ax}_T(x)$ | ‘ x is a logical or nonlogical axiom of T ’ |
| $\text{Prf}_T(x, y)$ | ‘ x is a proof of y in T ’ |

From the canonical proof predicate $\text{Prf}_T(x, y)$ one defines the provability predicate $\text{Pr}_T(x)$ as $\exists y \text{Prf}_T(x, y)$, expressing that there is a proof (a sequence of axioms or formulas obtained from the axioms by applications of the rules of inference) of x in T . The Π_1 -sentence $\text{Con}(T)$, expressing an intensionally correct consistency statement for T , is then simply $\text{Con}(T) : \Leftrightarrow \neg \text{Pr}_T(\perp)$, where \perp codes a falsity in T .

In this paper we will mostly deal with the language \mathcal{L} of EA expanded with one (or more) truth predicates and with extensions T of EA with axioms governing the behaviour of these predicates. In formulating the truth axioms in a relational setting, it is useful to resort to the fact that EA can represent the language ‘ \mathcal{L} plus domain constants’ via an injection $x \mapsto c_x$ formally associating each object x in the sense of T with these new constants. We will still call \mathcal{L} this expanded formal language internally represented in \mathcal{L} itself and write \bar{x} to denote these formal objects; the new formal constants will enable us to quantify into Gödel corners in the usual way. We write $\ulcorner \varphi(\bar{x}) \urcorner$ or $\text{sub}(\ulcorner \varphi(v) \urcorner, \bar{x})$ for $\text{sub}(\ulcorner \varphi(v) \urcorner, \ulcorner v \urcorner, \bar{x})$.

A sound translation function $\tau : \mathcal{L}_P(= \mathcal{L} \cup \{P\}) \rightarrow \mathcal{L}_1$, applied to sentences of the form $P^\ulcorner Pt \urcorner$ and replacing $P(\cdot)$ with some \mathcal{L}_1 -formula $\xi(\cdot)$, should of course yield $\xi(\tau^\ulcorner Pt \urcorner)$ and not $\xi(\ulcorner Pt \urcorner)$, where the notation $\tau(\cdot)$ refers to a functional expression in \mathcal{L} representing τ in EA. To achieve the required translation, one may resort to the recursion theorem (see §11.2 of (Rogers 1987)) that applies equally well to elementary functions; it yields for any recursive $f(x, y)$ an index e such that $f(e, y) \cong \phi_e(y)$, where $\phi_{(\cdot)}(\cdot)$ is the universal program. We can then define an elementary translation function τ_0 such that, in the relevant case, $\tau_0(x, P^\ulcorner Pt \urcorner) = \xi(\phi_{\ulcorner x \urcorner}(\ulcorner Pt \urcorner))$ and apply the recursion theorem to find an index e for τ_0 such that $\phi_e(P^\ulcorner Pt \urcorner) = \xi(\phi_e(\ulcorner Pt \urcorner))$. We can then let $\tau(x) \cong \phi_e(x)$.

When not otherwise specified, EA will be the base theory for our theories of truth. I do not attempt at giving sufficient conditions for what counts for a theory of truth, although I will often employ this expression. I hope that my choices will be self-explanatory and at any rate they are based on widely accepted choices.

I will employ different measures of the complexity of formulas: the *positive complexity* $|\varphi|^+$ of a formula φ is 0 for atomic and negated atomic formulas; it is $|\psi|^+ + 1$ if φ is $\neg\neg\psi$ or $\forall x\psi$, it is $\max(|\chi|^+, |\psi|^+) + 1$ if φ is $\chi \wedge \psi$. The *logical complexity* of a formula φ – formalized as $\text{lc}(\ulcorner \varphi \urcorner)$ – is simply the number of its logical symbols. These measures of complexity correspond to elementary functions and can be naturally represented in EA.

2.2. Interpretations and Isomorphisms The notion of relative interpretation and its variants will be central in this work, so they deserve careful introduction. We will consider only relational languages with finite signatures.

A *relative translation* of \mathcal{L}_T into \mathcal{L}_W can be described as a pair (δ, F) where δ is a \mathcal{L}_W -formula with one free variable – the domain of the translation – and F is a (finite) mapping that takes n -ary relation symbols of \mathcal{L}_T and gives back formulas of \mathcal{L}_W with n

free variables. The translation extends, modulo suitable renaming of bound variables,¹⁶ to the mapping τ :

- $(R(x_1, \dots, x_n))^\tau : \leftrightarrow F(R)(x_1, \dots, x_n)$;
- τ commutes with propositional connectives;
- $(\forall x \varphi(x))^\tau : \leftrightarrow \forall x (\delta(x) \rightarrow \varphi^\tau)$.

DEFINITION 2.1. *An (one-dimensional) interpretation K is specified by a triple (T, τ, W) , where τ is a translation of \mathcal{L}_T in \mathcal{L}_W , such that for all formulas $\varphi(x_1, \dots, x_n)$ of \mathcal{L}_T with the free variables displayed, we have:*

$$\text{if } T \vdash \varphi(x_1, \dots, x_n), \text{ then } W \vdash \bigwedge_{i=1}^n \delta_K(x_i) \rightarrow \varphi^\tau$$

I write $K : T \rightarrow W$ for ‘ K is an interpretation of T in W ’. An interpretation is *direct* if it maps identity to identity and it does not relativize quantifiers. I will often not distinguish between an interpretation and the relative translation that supports it. T and W are said to be *mutually interpretable* if there are interpretations $K : T \rightarrow W$ and $L : W \rightarrow T$. Since interpretability is a partial preorder on theories, mutual interpretability is an equivalence relation (its equivalence classes are known as degrees of interpretability).

I will often refer to a useful, model-theoretic characterization of interpretability. If $K : T \rightarrow W$ and \mathcal{M} is any model of W , then K can be seen as a method for defining a model \mathcal{M}^K of T inside \mathcal{M} .

LEMMA 2.2. *If $K : T \rightarrow W$ is identity preserving and W has full induction, then for any $\mathcal{M} \models W$ we find a (uniformly) \mathcal{M} -definable embedding of \mathcal{M} into an initial segment of \mathcal{M}^K .¹⁷*

Lemma 2.2. is obtained by noticing that in \mathcal{M} we can define an injection $f : \mathcal{M} \rightarrow \mathcal{M}^K$ such that, again by employing a convenient functional notation,

$$\begin{aligned} f(0^\mathcal{M}) &\mapsto 0^{\mathcal{M}^K}; & f(x +^\mathcal{M} 1^\mathcal{M}) &\mapsto f(x) +^{\mathcal{M}^K} 1^{\mathcal{M}^K} \\ f(x +^\mathcal{M} y) &\mapsto f(x) +^{\mathcal{M}^K} f(y) & f(x \times^\mathcal{M} y) &\mapsto f(x) \times^{\mathcal{M}^K} f(y) \end{aligned}$$

Full induction in \mathcal{M} is needed to prove the totality of f . f is clearly an isomorphism of \mathcal{M} in \mathcal{M}^K .

Given an intensionally correct consistency statement $\text{Con}(T)$, for T elementary presented,¹⁸ one can construct a formalized Henkin model of T in $\text{EA} + \text{Con}(T)$: the model is constructed in the standard way by means of Henkin axioms and it is equipped with a truth predicate – though not a truth predicate working for *all sentences* in the sense of T :

LEMMA 2.3. *Let the axiom set of T be captured by an elementary predicate. Then, $\text{EA} + \text{Con}(T)$ interprets T .*

To introduce stronger notions of interpretations, I introduce compositions of interpretations. Given $\tau_0 : \mathcal{L}_T \rightarrow \mathcal{L}_W$ and $\tau_1 : \mathcal{L}_W \rightarrow \mathcal{L}_V$, the composite of $K = (T, \tau_0, W)$ and $L = (W, \tau_1, V)$ is the interpretation $L \circ K = (T, \tau_1 \circ \tau_0, V)$, where $\delta_{L \circ K}(x) : \leftrightarrow \delta_K^L(x) \wedge \delta_L(x)$. Two

¹⁶ For more details on how clashes are avoided, I refer to (Visser 1997).

¹⁷ The condition on identity is redundant although not obviously so.

¹⁸ The same holds even if we allow T to be presented by a p-time set of axioms. A proof by Albert Visser of this stronger claim can be found in (Visser 1991). See also (Nicolai 2016).

interpretations $K_0, K_1 : T \rightarrow W$ are *equal* if W , the target theory, proves this. In particular, one requires,

$$\begin{aligned} W \vdash \forall x (\delta_{K_0}(x) \leftrightarrow \delta_{K_1}(x)) \\ W \vdash \forall \vec{x} (R_{K_0}(\vec{x}) \leftrightarrow R_{K_1}(\vec{x})) \quad \text{for any relation symbol } R \text{ of } \mathcal{L}_T \end{aligned}$$

A W -definable *morphism* between (again, one-dimensional) interpretations $K_0, K_1 : T \rightarrow W$ is a triple (K_0, ϕ, K_1) , with ϕ a \mathcal{L}_W -formula with two free variables, satisfying:

- (2) $W \vdash \forall x, y (\phi(x, y) \rightarrow (\delta_{K_0}(x) \wedge \delta_{K_1}(y)))$
- (3) $W \vdash \forall x, y, u, v (x =_{K_0} y \wedge u =_{K_1} v \wedge \phi(y, u) \rightarrow \phi(x, v))$
- (4) $W \vdash \forall x (\delta_{K_0}(x) \rightarrow \exists y (\delta_{K_1}(y) \wedge \phi(x, y)))$
- (5) $W \vdash \forall x, y, z (\phi(x, y) \wedge \phi(x, z) \rightarrow y =_{K_1} z)$
- (6) $W \vdash \forall x_1, \dots, x_n \forall y_1, \dots, y_n \left(\bigwedge_{i=1}^n \phi(x_i, y_i) \wedge R_{K_0}(x_1, \dots, x_n) \rightarrow R_{K_1}(y_1, \dots, y_n) \right)$

for any relation $R \in \mathcal{L}_T$.

To obtain an *isomorphism* from K_0 to K_1 one needs to add the following, extra conditions:

- (7) $W \vdash \forall y (\delta_{K_1}(y) \rightarrow \exists x (\delta_{K_0}(x) \wedge \phi(x, y)))$
- (8) $W \vdash \forall x, y, z (\phi(x, y) \wedge \phi(z, y) \rightarrow x =_{K_0} z)$
- (9) $W \vdash \forall x_1, \dots, x_n \forall y_1, \dots, y_n \left(\bigwedge_{i=1}^n \phi(x_i, y_i) \wedge R_{K_1}(y_1, \dots, y_n) \rightarrow R_{K_0}(x_1, \dots, x_n) \right)$

for any relation $R \in \mathcal{L}_T$.

We write $F : K_0 \cong K_1$ for ‘ F is an isomorphism from the interpretation K_0 to K_1 ’.

2.3. Sameness of Theories In this subsection I introduce two noticeable notions of sameness of theories extending extensional identity: bi-interpretability and synonymy (cf. (Visser 2006), (Friedman & Visser 2014)). They will play an important role in what follows.

DEFINITION 2.4. (SYNONYMY) *U and V are synonymous if and only if there are interpretations $K : U \rightarrow V$ and $L : V \rightarrow U$ such that V proves that $K \circ L$ and id_V are equal and, symmetrically, U proves that $L \circ K$ is equal to id_U .*

One can check that for U and V formulated in disjoint signatures, U and V are synonymous if and only if they have a common definitional extension. As shown in (Kaye & Wong 2007), for instance, ZF minus the axiom of infinity plus its negation and ‘every set has a transitive closure’ is synonymous with PA.¹⁹

To introduce a slightly less strict notion of equivalence, *bi-interpretability*, we first consider the notion of a retract, that in turn relies on the notion of isomorphism of interpretations introduced in the previous section.

DEFINITION 2.5. (RETRACT) *U is a retract of V if and only if there are $K : U \rightarrow V$ and $L : V \rightarrow U$ and a U -definable isomorphism between $L \circ K$ and id_U .*

¹⁹ ZF minus infinity plus its negation, rather surprisingly, is neither synonymous nor bi-interpretable with PA, (cf. (Enayat et alii 2010)).

The category-theoretic lexicon is due to Visser’s systematization of various notions of equivalence of theories in terms of appropriate categories of theories and interpretations (Visser 2006). We refer to Visser’s outstanding paper for further insights on the category-theoretic presentation. The notion of retract can also be visualized in model-theoretic terms: as it is illustrated in Figure 1, when $K: U \rightarrow V$ and $L: V \rightarrow U$, U is a retract of V if in any model $\mathcal{M} \models U$ we can construct an internal $\mathcal{M}^L \models V$ which in turn defines a model $(\mathcal{M}^L)^K \models U$ and there is an \mathcal{M} -definable isomorphism $F: (\mathcal{M}^L)^K \cong \mathcal{M}$.

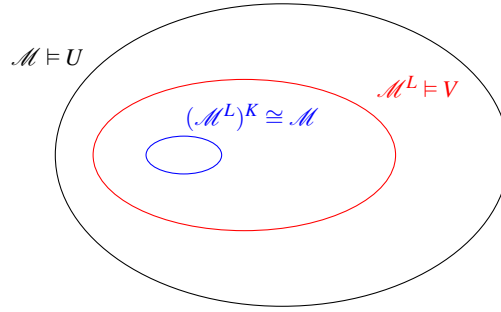


Fig. 1. U is a retract of V .

Essentially, U and V are bi-interpretable if and only if U is a retract of V and V is a retract of U using the same pair of interpretations.

DEFINITION 2.6. (BI-INTERPRETABILITY) *Given a pair of interpretations $K: U \rightarrow V$ and $L: V \rightarrow U$, U and V are bi-interpretable if and only if (i) there is a \mathcal{L}_V -formula F_0 such that V proves F_0 to be an isomorphism between $K \circ L$ and id_V and (ii) there is an \mathcal{L}_U -formula F_1 such that U proves F_1 to be an isomorphism between $L \circ K$ and id_U .*

Clearly, if U and V are synonymous, they are also bi-interpretable. The converse is, provably, not true. In Friedman & Visser (Friedman & Visser 2014) it is shown, for instance, that Adjunctive Set theory and its two-sorted, flattened version with a Frege function are bi-interpretable but not synonymous. Bi-interpretability preserves many mathematical properties of the theories: decidability, finite axiomatizability (cf. Lemma 2.8.), κ -categoricity. That synonymy is a stricter notion than bi-interpretability follows from the fundamental insight that the latter does preserve the automorphism group of models modulo isomorphism but not the action of the automorphism group on the domain of the model. Synonymy preserves both.

An important fact linking synonymy and bi-interpretability is the following, involving the so-called *sequential* theories.²⁰

LEMMA 2.7. ((Friedman & Visser 2014)) *Let U, V be sequential. If $K: U \rightarrow V$ and $L: V \rightarrow U$ witness a bi-interpretation of U and V and L is unrelativized and identity preserving, then U and V are synonymous.*

As hinted above, synonymy and interpretability can be seen, respectively, as equality and isomorphism in appropriate categories INT_0 and INT_1 of theories and interpretations (cf. (Visser 2006)). A category-theoretic framework is arguably the best way to formulate a

²⁰ A theory is sequential if it directly interprets Adjunctive set theory. More informally, a theory is sequential if it has a good coding of sequences. See (Friedman & Visser 2014) for details.

general theory of inter-theoretic reductions: in this paper, however, we will only employ the category-theoretic jargon as convenient notation. Since it will be relevant in what follows, we conclude the subsection with a sketch of the proof that bi-interpretability (and thus, synonymy), preserves finite axiomatizability.

LEMMA 2.8. (Visser 2006) *Let U, V be theories in finite signatures. Assume that $K: U \rightarrow V$ and $L: V \rightarrow U$ are interpretations and that U defines an isomorphism F from $L \circ K$ to id_U . Assume further that V is finitely axiomatizable. Then U is finitely axiomatizable.*

Proof Sketch. Let V_0 be the conjunction of a finite axiomatization of V . A finite $U_0 \subseteq U$ is specified by the single sentences: (i) F is an isomorphism between $L \circ K$ and id_U ; (ii) V_0^L . The theory U_0 is, by definition of retract, a subtheory of U . For the converse direction, one verifies that if U proves the sentence φ , then $U_0 \vdash \varphi^{K^L}$ by (ii) and the definition of retract. Thus $U_0 \vdash \varphi$ by (i). \square

§3. The usual suspects One option to formally capture the notion of conceptual equivalence of truth predicates may be to consider the *truth-free consequences* of the theories that characterize them. This is notoriously a bad idea. For my purposes here, it is useful to support this claim by introducing two theories that will be employed also in later sections.

DEFINITION 3.9. *The theory EA^\top in $\mathcal{L}_\top := \mathcal{L} \cup \{\top\}$ is obtained by formulating EA in the new language and allowing the truth predicate into instances of Δ_0 -Ind. EA^\top , in turn, only features arithmetical Δ_0 -induction.*

DEFINITION 3.10. (TB) *The theory TB is obtained by extending EA^\top with all instances of the schema*

$$(tb) \quad \top^\top \varphi^\top \leftrightarrow \varphi$$

for φ a sentence of \mathcal{L} . TB^\dagger is obtained by adding (tb) to EA^\top .

The other theory is essentially the axiom set introduced by Feferman and inspired to the Strong Kleene version of Kripke's theory of truth (Cantini 1989; Feferman 1991).

DEFINITION 3.11. (KF \dagger) *KF \dagger extends EA^\top with the universal closures of the following:*

- Tat $(\top^\top R(\dot{x}_1, \dots, \dot{x}_n)^\top \leftrightarrow R(x_1, \dots, x_n)) \wedge (\top^\top \text{Tng}(\top^\top R(\dot{x}_1, \dots, \dot{x}_n)^\top) \leftrightarrow \neg R(x_1, \dots, x_n))$
- T1 $\text{Sent}_{\mathcal{L}_\top}(x) \rightarrow (\top^\top \text{dn}(x) \leftrightarrow \top x)$
- T2 $\text{Sent}_{\mathcal{L}_\top}(\text{and}(x, y)) \rightarrow (\top^\top \text{and}(x, y) \leftrightarrow \top x \wedge \top y)$
- T3 $\text{Sent}_{\mathcal{L}_\top}(\text{and}(x, y)) \rightarrow (\top^\top \text{ng}(\text{and}(x, y)) \leftrightarrow \top^\top \text{ng}(x) \vee \top^\top \text{ng}(y))$
- T4 $\text{Sent}_{\mathcal{L}_\top}(\text{all}(v, x)) \rightarrow (\top^\top \text{all}(v, x) \leftrightarrow \forall y \top^\top \text{sub}(x, \bar{y}))$
- T5 $\text{Sent}_{\mathcal{L}_\top}(\text{all}(v, x)) \rightarrow (\top^\top \text{ng}(\text{all}(v, x)) \leftrightarrow \exists y \top^\top \text{sub}(\text{ng}(x), \bar{y}))$
- T6 $\top^\top \top \dot{x}^\top \leftrightarrow \top x$
- T7 $\top^\top \text{Tng}(\top \dot{x}^\top) \leftrightarrow \top^\top \text{ng}(x)$
- T8 $\top x \rightarrow \text{Sent}_{\mathcal{L}_\top}(x)$

The arguments that witness the conservativeness of TB^\dagger and KF^\dagger over EA are well-known. For TB^\dagger one replaces in every derivation in TB^\dagger of a truth-free formula the truth predicate with an EA-definable, partial truth predicate. In the second case, starting with any model of EA (in the language with domain constants), we expand it via a positive inductive definition of a Kripke truth set to a model of KF^\dagger (see §4 of (Cantini 1989)).

LEMMA 3.12. $TB\upharpoonright$, EA , and $KF\upharpoonright$ prove the same \mathcal{L} -theorems.

By contrast, it is not hard to see that $KF\upharpoonright$ and $TB\upharpoonright$ fare much differently with respect to the conditions in §2. $TB\upharpoonright$ is typed and $KF\upharpoonright$ is type-free; $KF\upharpoonright$ is compositional to a large extent – only negation, for obvious reasons, does not fully commute with T – whereas $TB\upharpoonright$ can hardly be thought as compositional: for instance, $TB\upharpoonright$ cannot prove the sentence

$$(10) \quad \forall x, y (\text{Sent}_{\mathcal{L}}(\text{and}(x, y)) \rightarrow (T(\text{and}(x, y)) \leftrightarrow (Tx \wedge Ty)))$$

The unprovability of (10) also suggests that TB falls short of many generalizations provable in $KF\upharpoonright$, suggesting that also this desideratum separates the two theories. Finally, TB and $KF\upharpoonright$ are theories whose truth predicates capture, in the sense of (Fischer et alii 2015) instantiated by (1), different semantic constructions: the Tarskian truth set and the fixed points of Kripke’s theory of truth respectively.

The case just considered is not isolated: It is not hard to find many other counterexamples to the claim that the provability of the same truth-free consequences amounts to an adequate formal rendering of the conceptual equivalence of two theories’ truth predicates.

An alternative may be represented by *mutual interpretability*. A little reflection shows, however, that mutual interpretability suffers from a problem that is in a sense opposite to the one suffered by provability of the same base-theoretic consequences: it mixes up syntactic and semantic aspects of our theories.

LEMMA 3.13. (Feferman) Any $T \supseteq EA$ interprets $T + \neg \text{Con}(T)$.

Proof Sketch. This proof is due to Visser and, independently, to Lindström; it does not rely on the reflexivity of the theory involved and therefore suits our purposes. By Gödel’s second incompleteness theorem, $T \vdash \text{Con}(T) \rightarrow \text{Con}(T + \neg \text{Con}(T))$, therefore $T + \text{Con}(T)$ interprets $T + \neg \text{Con}(T)$ by Lemma 2.3., let’s say with an interpretation K . Finally we combine the identity interpretation on $T + \neg \text{Con}(T)$ and K : if $\text{Con}(T)$, pick K , if $\neg \text{Con}(T)$, $\text{id}_{T + \neg \text{Con}(T)}$ suffices. \square

By Lemma 3.13., any theory of truth over EA will mutually interpret its inconsistency. Since we usually require syntactic notions – including $\text{Con}(T)$ – to be canonically constructed and T to be ω -consistent, any theory of truth will be ‘equivalent’ with a theory that compromises our basic assumptions on the structure of syntactic universe.²¹ Mutual interpretability does not guarantee that the structure of the so-called bearers of truth remains fixed across theories. For the reasons explained in §2., this seems to undermine the very possibility of comparing truth predicates.²²

²¹ An advocate of mutual interpretability, at this point, may highlight the artificiality of the theories like $T + \neg \text{Con}(T)$. After all, they are artifacts just right to reach Σ_1 -unsoundness, therefore ω -inconsistency. However, more natural examples are not difficult to construct. To mention a familiar example, let us consider the theory $TB[PA]$ just introduced and the theory resulting from extending PA with restricted induction with Tat , the typed versions of $T2$, $T4$ and

$$(11) \quad \text{Sent}_{\mathcal{L}}(x) \rightarrow (Tng(x) \leftrightarrow \neg Tx),$$

known in the literature as $CT\upharpoonright[PA]$ or PA^{FT} . $TB[PA]$ and $CT\upharpoonright[PA]$ are mutually interpretable, as $CT\upharpoonright[PA]$ is interpretable in PA (cf. (Enayat & Visser 2015)), but the theories fare rather differently in terms of the conditions in §2. and do not share the same conception of truth in the sense defined above.

²² (Fujimoto 2010) already employed Feferman’s theorem on the interpretability of inconsistency – applied to TB – to argue against mutual interpretability as a notion of theoretical equivalence for theories of truth. For a recent, thorough study of Feferman’s theorem, see Visser (Visser 2015).

To overcome these and further difficulties with assuming provability of the same base-theoretic consequences or mutual interpretability as formal counterparts of conceptual equivalence of two theories of truth, we may resort to the notion of *relative truth-definability* due to (Fujimoto 2010) and (Halbach 2014). In a nutshell, T is truth-definable in W if there is a direct, relative interpretation of T in W that preserves the vocabulary of the syntactic base theory B . I give a more pedantic definition to fix some notation, starting with an arbitrary base theory B containing EA. It is convenient here to consider the language of our theories of truth as extending the base language \mathcal{L}_B with an indexed set of semantic predicates $\{S_i\}_{i \in I}$: we want to allow in fact for the possibility that truth is expressed via, for instance, a binary satisfaction predicate.

DEFINITION 3.14. (TRUTH TRANSLATION) A truth-translation $\tau: \mathcal{L}_T \rightarrow \mathcal{L}_W$, where $\mathcal{L}_T = \mathcal{L}_B \cup \{S_i\}_{i \in I}$ and $\mathcal{L} \in \mathcal{L}_T \cap \mathcal{L}_W$, is a pair (d, F) , where $d(x)$ is $x = x$ and F a mapping of n -ary predicates of \mathcal{L}_T into \mathcal{L}_W formulas with n free variables satisfying:

$$\begin{aligned} \tau(R)(\vec{x}) &: \leftrightarrow R(\vec{x}) && \text{for } n\text{-ary } R \in \mathcal{L}_B; \\ \tau(S_i)(\vec{x}) &: \leftrightarrow F(S_i)(\vec{x}) && \text{for } i \in I \\ \tau(\neg\varphi) &: \leftrightarrow \neg\tau(\varphi) \\ \tau(\varphi \wedge \psi) &: \leftrightarrow \tau(\varphi) \wedge \tau(\psi) \\ \tau(\forall x\varphi) &: \leftrightarrow \forall x(d(x) \rightarrow \tau(\varphi)) \end{aligned}$$

The mapping F behaves like the identity function when applied to base-theoretic, nonlogical symbols. Moreover, the translation is clearly unrelativized.

DEFINITION 3.15. (RELATIVE TRUTH-DEFINITION) A truth-definition of a theory T in $\mathcal{L}_T = \mathcal{L}_B \cup \{S_i\}_{i \in I}$ – with I an index set – in a theory W is a triple (T, τ, W) , where τ is a truth-translation of \mathcal{L}_T into \mathcal{L}_W and for all \mathcal{L}_T -sentences φ , if $T \vdash \varphi$, then $W \vdash \tau(\varphi)$.

Like relative interpretability, relative truth-definability is a partial preorder. We call two theories of truth T and W *mutually truth-definable* if there are truth-definitions $K: T \rightarrow W$ and $L: W \rightarrow T$. Mutual truth-definability defines indeed an equivalence relation on theories. As (Fujimoto 2010) points out, it is immediate from the definition that if W truth-defines T and any $\mathcal{M} \models B$ admits an expansion to a model of W , then \mathcal{M} admits an expansion to a model of T .²³ Also, due to the non-relativization of quantifiers, if T is Σ_1 -unsound or ω -inconsistent and T is truth-definable in W , W will be Σ_1 -unsound or ω -inconsistent. Therefore, mutual truth-definability prevents the possibility of equating theories with non-isomorphic ontological commitments. Also, it takes care of what the theory proves true as if T truth defines W , there should be a formula of the language of T defining the extension of the truth predicate of W – in the case of a single truth predicate in W – simulating the behaviour of this truth predicate. But is this enough to conclude that the truth predicates of mutually truth-definable theories embody the same conception of truth? There are well-known examples, it seems, that suggest a negative answer to this question. We consider some examples.

In the following a formula of \mathcal{L}_T will be said to be *T-positive* if and only if T does not occur, in φ , in the scope of an odd number of negation signs. (Halbach 2009) introduced

²³ This is simply because the numbers of the two models are the same, and the truth-definition ensures that the right extension of the truth predicates of T can be always extracted via the extension of the truth predicate of W .

and motivated a collection of type-free, uniform T -sentences that escape the Liar paradox by controlling the behaviour of negation.

DEFINITION 3.16. PUTB is the theory obtained by extending EA^T with the schema

$$(putb) \quad \forall x (T(\ulcorner \varphi(x) \urcorner) \leftrightarrow \varphi(x))$$

for all T -positive formulas φ of \mathcal{L}_T .

DEFINITION 3.17. KF is the theory obtained by extending EA^T with Tat-T7.

LEMMA 3.18. PUTB is a subtheory of KF.

Lemma 3.18. is obtained by showing, via a meta-induction on the positive complexity of \mathcal{L}_T -formulas, that (putb) is provable in KF. Trivially, therefore, KF defines the truth predicate of PUTB. More interestingly, (Halbach 2009) also shows that

LEMMA 3.19. PUTB defines the truth predicate of KF.

Proof Sketch. EA^T proves the diagonal lemma in parametrized form, i.e.: for any \mathcal{L}_T -formula $\varphi(x, \vec{y})$, we find a unary formula $\psi(\vec{y})$ such that

$$(12) \quad EA^T \vdash \varphi(\vec{y}) \leftrightarrow \psi(\ulcorner \varphi(\vec{y}) \urcorner, \vec{y})$$

Now if we let $\mathcal{I}(x, Ty)$ be a \mathcal{L}_T -formula that mimics the positive inductive definition of a Kripke truth-set under the scope of T (see p. 190 of (Halbach 2014)), (12) will yield a formula $\theta(x)$ such that

$$(13) \quad EA^T \vdash \theta(x) \leftrightarrow \mathcal{I}(x, T\ulcorner \theta(x) \urcorner)$$

Since in $\theta(x)$ there are only positive occurrences of T , they can be removed so that one can verify that $\theta(x)$ satisfies all axioms of KF. \square

If mutual truth-definability were a sufficient condition to establish the conceptual equivalence of the theories' truth predicates, the concepts of truth associated with PUTB and KF would be equivalent. This is the position taken, for instance, by (Fujimoto 2010).²⁴ Despite their mutual truth-definability, PUTB and KF fare much differently with respect to the criteria highlighted in §2.. Firstly, PUTB is can hardly be considered to be *compositional* and it cannot prove as many general claims as KF by using its own notion of truth. As shown in Lemma 6.1, in fact, any proof in PUTB only employs positive disquotational

²⁴ He writes:

...[the mutual truth-definability of KF and PUTB] may be still interpreted to indicate that KF and PUTB are indeed 'conceptually equivalent' despite superficial differences. (p. 342)

To be fair with Fujimoto's position, he also claims:

...we propose [relative truth-definability, A/N] by means of which we can represent a certain 'equivalence' or 'reducibility' among theories of truth from a more conceptual point of view, although it may be still too coarse to fully represent 'conceptual equivalence' or 'conceptual reducibility' (in some strong sense). (Fujimoto 2010)

The issue thus becomes dangerously verbal. At any rate, even if we take mutual truth-definability to witness a weak form of 'conceptual equivalence', what I will say in what follows may be seen as an invitation to investigate and an attempt to regiment stronger notions of conceptual equivalence.

sentences applied to \mathcal{L}_T -sentences φ with limited, standard complexity. But there is more: PUTB can be consistently made symmetric, that is the result of closing PUTB under the rules

$$(T\text{-in}) \quad \frac{\varphi(x)}{T \vdash \varphi(x)} \qquad (T\text{-out}) \quad \frac{T \vdash \varphi(x)}{\varphi(x)}$$

for $\varphi(v) \in \mathcal{L}_T$, yields a consistent theory.²⁵ It is in fact possible to rule out putative proofs of inconsistencies in PUTB plus (T-in) and (T-out) in the following way: if (T-out) and (T-in) are applied to T-positive sentences they can be dispensed with via (putb). Moreover, in any proof in PUTB plus (T-in) and (T-out) — where the two rules are applied to sentences that are not T-positive — (i) applications of (T-out) may be eliminated; (ii) the semantics of PUTB can be adapted to accommodate possible applications of (T-in) (see (Halbach 2014), §19.5, for proofs of these facts). By contrast, KF is essentially asymmetric, as it cannot be consistently closed under the two rules. Finally, if KF is sound with respect to fixed points of the four-valued version of Kripke’s semantic construction, PUTB is not as it has models of the form (\mathbb{N}, S) where S is not such a fixed point.

Mutual truth-definability, therefore, does not preserve compositionality (in the sense specified above), provable generalizations, symmetry, at least if one takes at face value *the* truth predicate(s) of the theories. Moreover, mutually truth-definable theories may not be equally close to a semantic construction. In the next section we will see that not even finite axiomatizability is preserved. Back to our example, there is no doubt that PUTB has the resources to replicate the behaviour of the truth predicate of KF via a non-primitive notion, but this would hardly count as evidence for the claim that PUTB and KF are conceptually equivalent. For this matter, only the primitive notion of truth of the theories should matter.

§4. Notions of Equivalence of Truth Predicates In this section I consider the varieties of equivalences of theories introduced in §2 and investigate whether they may help in overcoming the problems encountered above. My strategy will be suitably adapting those definitions to obtain notions that

- (I) contain mutual truth-definability – that is, the latter notion will play the role of necessary condition for the new notions;
- (II) properly extend mutual truth-definability so to capture the differences, for instance, between the truth predicates of PUTB and KF;
- (III) are non-empty, in the positive sense that there are natural theories of truth that fall into the relations.

Since this is in a sense unexplored territory, it is not my main focus to find notions of inter-theoretic reduction that are *just right* to capture conceptual equivalence of truth predicates: by contrast, I look for strict notions that may eventually be relaxed and calibrated to accommodate our philosophical purposes. In practice, I will proceed as follows: the notions of retract, bi-interpretability and synonymy from §2.2. are, by definition, not in continuity with mutual truth-definability. This is simply because they are not fine tuned for theories extending a common syntax theory with truth axioms: that is, they do not require the interpretations witnessing them to be unrelativized and behaving like the identity on the

²⁵ Of course, to avoid the Liar paradox, (T-in) and (T-out) need to be understood as rules applying only to *theorems* and not to assumptions.

syntactic vocabulary. There may be no such a thing as syntactic vocabulary in them. Therefore, I modify the definitions of retract and bi-interpretation accordingly by requiring the composed interpretations to be truth-definitions. I will also not forget about the notion of bi-interpretability simpliciter and consider examples of theories of truth that are bi-interpretable but not mutually truth-definable.

DEFINITION 4.20. (T-RETRACT) *Let T, W be theories of truth based on a syntactic base theory B . T is a t -retract of W if there are truth-definitions $K: T \rightarrow W$ and $L: W \rightarrow T$ and a T -definable isomorphism between $L \circ K$ and id_T .*

At the intuitive level, if T is a t -retract of W , then W functions as a faithful mirror for T : the image of T reflected via its truth-definition in W is a faithful one modulo isomorphism. There is no guarantee, however, that also the converse is true: it may well be the case that, when W looks for its image reflected in T via its truth-definition in W , the result is a distorted image. The notion of t -equivalence, by contrast, requires both theories to be ‘faithful mirrors’ focusing on the same pairs of truth-definitions:

DEFINITION 4.21. (T-EQUIVALENCE) *Let T, W be theories of truth based on B . T and W are truth equivalent iff there are truth-definitions $K: T \rightarrow W$ and $L: W \rightarrow T$ such that (i) T proves that there is an isomorphism between $L \circ K$ and id_T ; (ii) W proves that there is an isomorphism between $K \circ L$ and id_W .*

It is immediate from the definitions that t -retractions and t -equivalences simply impose further conditions on mutual truth-definitions. Focusing on them, the fulfillment of desideratum (I) on page 16 is then immediately obtained.

LEMMA 4.22. *If T is a t -retract of W , then T and W are mutually truth-definable. A fortiori, if T and W are t -equivalent, then they are mutually truth-definable.*

Lemma 2.7. tells us that, since we are dealing with direct interpretations, t -equivalence is a form of synonymy. I therefore propose the following, quasi-formal thesis:

THESES 1 *If two theories are t -equivalent, their associated semantic notions are conceptually equivalent.*

Thesis 1 only imposes sufficient conditions on the conceptual equivalence of semantic notions. Obviously, it is not a formally precise claim, although its plausibility seems to be safely grounded in the strictness of the notion of synonymy. Nonetheless, it requires a suitable amount of empirical data to be confirmed or dismissed: in the following sections I start collecting these data by showing that Thesis 1 captures the non-conceptual equivalence of KF and PUTB and that there are simple, reassuring cases of t -equivalence.

4.1. Separating mutual truth-definability, t -retracts, t -equivalence. In this section I consider desideratum (II) on page 16 and show that the notions of t -retract and t -equivalence are properly stricter than mutual truth-definability in a formally precise sense.

PROPOSITION 4.23. *PUTB is not finitely axiomatizable.*

Proof. Let’s assume that there is a sentence A such that $A \dashv\vdash \text{PUTB}$. By the finite axiomatizability of EA, A has the normal form $A_0 + A_1$, where A_0 is a finite set of instances of (putb) and A_1 is a finite version of EA^T . The instances of (putb) in A_0 are therefore applied only to sentences ϕ such that $\text{lc}(\ulcorner \phi \urcorner) \leq n$ (the number of logical symbols in ϕ), for a standard n . By adapting the semantics of PUTB given in (Halbach 2009) I show that there is a model of A that is not a model of PUTB.

Let $\text{Sent}_{\mathcal{L}_T^+}^n$ be the set of T-positive sentences of \mathcal{L}_T with no more than n logical symbols (excluding the identity symbol), for a suitable $n \in \omega$ and let

$$C : \leftrightarrow \neg \underbrace{(0 = 0 \wedge \dots \wedge 0 = 0)}_{\wedge \text{ applied } n\text{-times}}$$

Notice that, despite the presence of the negation symbol, C is trivially T-positive. We define the operator $\Phi: \mathcal{P}(\omega) \rightarrow \mathcal{P}(\omega)$:

$$\Phi(S_1) := \{\varphi \in \text{Sent}_{\mathcal{L}_T^+}^n \mid (\mathbb{N}, S_1) \models \varphi\} \cup \{m \in S_1 : m \notin \text{Sent}_{\mathcal{L}_T^+}^n\} \cup \{C\}$$

It should be clear that $\Phi(\cdot)$ is monotone. The monotonicity of Φ entails the existence of fixed points of Φ and, in particular, of its minimal fixed point \mathcal{I}_Φ obtained by closing the empty set under Φ . It should be noticed that $C \in \mathcal{I}_\Phi$ (and it is in any fixed point of Φ). Moreover, $(\mathbb{N}, \mathcal{I}_\Phi)$ is a model of A . The axioms of EA clearly hold in $(\mathbb{N}, \mathcal{I}_\Phi)$. For $\varphi \in \text{Sent}_{\mathcal{L}_T^+}^n$, we have

$$\begin{aligned} (\mathbb{N}, \mathcal{I}_\Phi) \models T\varphi &\Leftrightarrow \varphi \in \mathcal{I}_\Phi (= \Phi(\mathcal{I}_\Phi)) \\ &\Leftrightarrow (\mathbb{N}, \mathcal{I}_\Phi) \models \varphi. \end{aligned}$$

To complete the proof of Proposition 4.23.: given the assumption of the extensional identity of A and PUTB, $(\mathbb{N}, \mathcal{I}_\Phi)$ would also be a model of PUTB. But then $(\mathbb{N}, \mathcal{I}_\Phi) \models T \ulcorner C \urcorner$, thus $(\mathbb{N}, \mathcal{I}_\Phi) \models C$. That is, the desired contradiction. \square

COROLLARY 4.24. $\text{PUTB} \upharpoonright$ is not finitely axiomatizable.²⁶

Since, by inspection, the theory that we called KF is finitely axiomatizable (recall that we are constructing KF over EA^T):

COROLLARY 4.25. *Relative truth-definability does not preserve finite axiomatizability.*

But, by Lemma 2.8., finite axiomatizability is preserved over retractions:

COROLLARY 4.26.

- (i) PUTB is not a retract and therefore not a t -retract of KF;
- (ii) PUTB and KF are not bi-interpretable, therefore also not t -equivalent.

The argument would still go through if we replace EA^T with EA^t . Therefore, Corollary 4.26. holds if we replace PUTB and KF with $\text{PUTB} \upharpoonright$ and $\text{KF} \upharpoonright$ respectively.

A similar line of reasoning can be applied to other cases, less relevant for our discussion but still worth considering.

DEFINITION 4.27.

- (i) $\text{CT} \upharpoonright$ is the theory in \mathcal{L}_T that extends EA^t with Tat , the restrictions of T2, T4, T8 to sentences of \mathcal{L} , and the sentence

$$(\star) \quad \text{Sent}_{\mathcal{L}}(x) \rightarrow (Tng(x) \leftrightarrow \neg Tx)$$

²⁶ The method of Proposition 4.23. can be applied more generally to yield the non finite axiomatizability of typed T-sentences over an arbitrary base theory in predicate logic. For instance, axioms of the form

$$Tc_i \leftrightarrow \varphi_i$$

for all $i \in \omega$ are not finitely axiomatizable over predicate logic PRED alone.

(ii) $\text{UTB}\upharpoonright$ is the theory obtained by extending EA^\dagger with the schema

$$(\text{utb}) \quad \forall x (\top^\top \varphi(x)^\top \leftrightarrow \varphi(x))$$

for all formulas $\varphi(v)$ of \mathcal{L} with one free variable.

By the finite axiomatizability of EA , it follows that $\text{CT}\upharpoonright$ is finitely axiomatizable. A simplified version of the argument in Proposition 4.23. shows that $\text{UTB}\upharpoonright$ is not finitely axiomatizable.

COROLLARY 4.28. *$\text{UTB}\upharpoonright$ is not a t -retract of $\text{CT}\upharpoonright$. Therefore, $\text{CT}\upharpoonright$ and $\text{UTB}\upharpoonright$ are not t -equivalent.*

Unlike Corollary 4.26., the failure of t -equivalence between $\text{UTB}\upharpoonright$ and $\text{CT}\upharpoonright$ does not amount to a further example separating mutual truth-definability and the stricter notions we are considering. Although $\text{UTB}\upharpoonright$ is trivially truth-definable in $\text{CT}\upharpoonright$, the converse fails.²⁷

I conclude this subsection with few remarks on desideratum (II) in the light of the observations above. Corollary 4.26. relies on the finite axiomatizability of KF (and $\text{KF}\upharpoonright$) that, in our setting, is easily reached by the finite axiomatizability of EA . In §2. we included the finite axiomatizability of a theory as a property that should be preserved between theories whose truth predicates are considered to be conceptually equivalent. This is a delicate point and we should pause a little more on it: It would seem in fact that two theories like KF and $\text{KF}[\text{PA}]$ – that is $\text{Tat} - \text{T8}$ added to EA^\top or to its version with full induction – *should* be conceptually equivalent; yet this is ruled out as $\text{KF}[\text{PA}]$ is not finitely axiomatizable. Can this asymmetry threaten the entire project of a formal analysis of conceptual equivalence between truth theories?

This concern can be addressed by elaborating on what has been already anticipated in §2.. We are considering comparisons of axiomatic systems, and there is considerable evidence for the claim that theories of truth built on different base theories are in a sense incomparable. What we compare are combinations of truth-theoretic and syntactic principles, where a syntax theory is fixed and ‘universal’.²⁸ If we allowed for a comparison of theories of truth built on non-equiconsistent base theories, for instance, we would violate some of the very criteria on which truth predicates are in general evaluated. $\text{CT}\upharpoonright[\text{PA}]$, for instance, is interpretable in PA ; $\text{CT}\upharpoonright$, by contrast, is provably not interpretable in EA . If we let $\text{CT}\upharpoonright[\text{PA}]$ and $\text{CT}\upharpoonright$ share the same conception of truth in the strong sense we are after, we would also be led to consider as conceptually equivalent a truth predicate that can be relatively interpreted in the base theory, and a truth predicate that cannot. Therefore two theories sharing the same conception of truth may be regarded in considerably different ways by, for instance, advocates of the expressive irreducibility of the notion of truth.

After all, it seems, it’s a sensible choice to keep the syntax theory fixed. But there seem to be reasons to go even further and motivate the choice of EA as base theory for this particular project. Its finite axiomatizability harmonizes with the choice of our truth axioms in a better way than other choices. By adding a suitable infinite set of truth axioms to EA it may be the case – and it is in the case of TB , UTB and PUTB as defined above – that we obtain a non finitely axiomatizable extension of EA ; by adding a finite set of truth axioms to it,

²⁷ Otherwise the theory CT obtained from $\text{CT}\upharpoonright$ by allowing \mathcal{L}_\top -formulas into the schema of Δ_0 -induction will be interpretable in UTB , namely $\text{UTB}\upharpoonright$ with Δ_0 - \top induction, quod non as CT proves $\text{Con}(\text{EA})$.

²⁸ By the considerations below it may be possible to allow for bi-interpretable syntax theories, but surely non-equiconsistent base syntactic theories should be disregarded.

like in the case of KF and CT, we obtain a finitely axiomatizable extension of EA. If other choices of the base theories are made, e.g. PA, it is often the case that the size of our set of truth axioms does not matter for the finite axiomatizability of the resulting theory of truth. $TB[PA]$, for instance, is not finitely axiomatizable. Despite the finite axiomatizability of CT , however, $CT \upharpoonright [PA]$ is not finitely axiomatizable essentially because of the reflexivity of PA (cf. (Fischer 2009; Leigh 2015; Enayat & Visser 2015)).

4.2. *t-equivalence, t-retracts, bi-interpretability: some examples.* In this section I move to desideratum (III) and introduce theories of primitive truth that are *t-retracts* of others or that are *t-equivalent* to others. By definition, *t-equivalence* satisfies the criteria for Lemma 2.7. and leads straight to synonymy (or definitional equivalence under minimal assumptions). Furthermore, the cases of *t-equivalence* that I will consider are in fact cases of equality of truth-definitions. One may think of introducing a notion of *t-synonymy* that relates to *t-equivalence* like synonymy simpliciter does to bi-interpretability. Given Visser and Friedman's result and the canonical definition of truth translation, however, this would not yield any further insights so I stick with the more general notion of *t-equivalence*.

I consider again extensions of EA^T and extensions of EA with a binary satisfaction predicate. Let $tr: \mathcal{L}_{Sat} \rightarrow \mathcal{L}_T$ and $st: \mathcal{L}_T \rightarrow \mathcal{L}_{Sat}$, where $\mathcal{L}_{Sat} := \mathcal{L} \cup \{Sat\}$ and S is a binary satisfaction predicate, be truth-translations specified by

$$\begin{aligned} Sat^{tr}(x, y) &:\leftrightarrow Tsb(x, y) \\ T^{st}(x) &:\leftrightarrow Sat(\langle \rangle, x) \end{aligned}$$

where $\langle \rangle$ denotes the empty sequence and the elementary formula $sb(x, y) = z$ captures the operation of substituting in a formula y its free variables with the result of applying the finite mapping x to the free variables of y . We assume the functions sub and sb to be defined in such a way that in EA we can prove:

$$(14) \quad sub(\ulcorner \varphi(v) \urcorner, \ulcorner v \urcorner, \bar{x}) = sb(\langle x \rangle, \ulcorner \varphi(v) \urcorner)$$

$$(15) \quad \forall x (sb(\langle \rangle, x) = x)$$

In investigating positive results on *t-retracts* and *t-equivalence*, I start with a simple example to fix the basic reasoning involved and then move to slightly more complex cases of the same sort. Let TBS be EA formulated in \mathcal{L}_{Sat} plus all instances of the schema

$$(tbs) \quad Sat(\langle \rangle, \ulcorner \varphi \urcorner) \leftrightarrow \varphi$$

for all \mathcal{L} -sentences φ .²⁹

PROPOSITION 4.29.

- (i) TB and TBS are mutually truth-definable;
- (ii) TB is a *t-retract* of TBS.

Proof. The only thing to check is (ii). Since we have truth-definitions, we need to take care only of the truth predicate. In TB, $T^{tr \circ st}x$ is equivalent to $Tsb(\langle \rangle, x)$ that, by (15), is equivalent to Tx . The required isomorphism is the identity on TB. \square

Here I am not able to employ tr and st to show that TBS is a retract of TB. The reason is that I cannot prove in TBS the expected interaction of substitution and the satisfaction predicate: informally, I cannot prove that the result of substituting the elements of (x_1, \dots, x_n)

²⁹ Here I consider the case in which Δ_0 -induction is extended to \mathcal{L}_{Sat} , but the same argument goes through for the corresponding theories with restricted induction.

in $\varphi(v_1, \dots, v_n)$ is satisfied by the empty sequence exactly when $\varphi(v_1, \dots, v_n)$ is satisfied by (x_1, \dots, x_n) . So I let

$$\text{TBS}^* := \text{TBS} \cup \{\forall x, y (\text{Sat}(\langle \rangle, \text{sb}(x, y)) \leftrightarrow \text{Sat}(x, y))\}$$

I call D the sentence in curly brackets. That TBS^* is a natural extension of TBS is justified by the following claim:

LEMMA 4.30. *TBS* and TB are mutually truth-definable.*

Proof. That TBS^* defines the truth predicate of TB follows from Proposition 4.29.. Also:

$$\begin{aligned} \text{TB} \vdash [\text{Sat}(\langle \rangle, \text{sb}(x, y)) \leftrightarrow \text{Sat}(x, y)]^{\text{tr}} &\leftrightarrow (\text{Tsb}(\langle \rangle, \text{sb}(x, y)) \leftrightarrow \text{Tsb}(x, y)) \\ &\leftrightarrow (\text{Tsb}(x, y) \leftrightarrow \text{Tsb}(x, y)) \end{aligned}$$

□

Therefore, we have our first, simple-minded example of t-equivalent theories of truth.

PROPOSITION 4.31.

- (i) TBS^* is a t-retract of TB;
- (ii) TB and TBS^* are t-equivalent.

Proof. In TBS, $\text{Sat}^{\text{stotr}}(x, y)$ is provably equivalent to $\text{Sat}(\langle \rangle, \text{sb}(x, y))$. In TBS^* we can then conclude $\text{Sat}(x, y)$. □

The role of D becomes even more prominent if we consider a suitably simplified form of uniform, disquotational satisfaction. UTB has already been introduced on page 19.

DEFINITION 4.32. *UTBS extends EA^T with all instances of the schema*

$$(\text{utbs}) \quad \forall x (\text{Sat}(\langle x \rangle, \ulcorner \varphi(v) \urcorner) \leftrightarrow \varphi(x))$$

for all \mathcal{L} -formulas φ with the displayed variables free. Similarly as before, we let

$$\text{UTBS}^* := \text{UTBS} + \text{D}.$$

By adapting the reasoning employed for TB and TBS, and by (14) above, we can conclude:

LEMMA 4.33.

- (i) UTBS^* and UTB are mutually truth-definable.
- (ii) UTBS^* and UTB are t-equivalent.

I now consider the case of full Tarskian truth and satisfaction over EA. In EA there are the following elementary formulas expressing the corresponding elementary syntactic notions:

$$\begin{aligned} \text{map}(x) &:\Leftrightarrow \text{'x is a finite mapping'} \\ \text{dom}(x, y) &:\Leftrightarrow \text{'y is an element of the domain of the finite mapping x'} \\ \text{ass}(x, y) &:\Leftrightarrow \text{'y is a formula and x is a finite mapping whose domain} \\ &\quad \text{contains precisely the free variables of y'} \\ x \supseteq y &:\Leftrightarrow \text{'x and y are assignments and dom}(x) \text{ contains dom}(y)' \end{aligned}$$

DEFINITION 4.34. *The theory CS in \mathcal{L}_{Sat} is obtained by extending EA formulated in \mathcal{L}_{Sat} – that is, in which Sat can appear into instances of induction – with the universal closures of the following:*

(S1)

$$\text{ass}(x, \ulcorner R(v_1, \dots, v_n) \urcorner) \wedge \left(\bigwedge_{i=1}^n x(\ulcorner v_i \urcorner) = x_i \right) \rightarrow (\text{Sat}(x, \ulcorner R(v_1, \dots, v_n) \urcorner) \leftrightarrow R(x_1, \dots, x_n))$$

(S2)

$$\text{ass}(x, y) \rightarrow (\text{Sat}(x, \text{ng}(y)) \leftrightarrow \neg \text{Sat}(x, y))$$

(S3)

$$\text{ass}(x, \text{and}(y, z)) \rightarrow (\text{Sat}(x, \text{and}(y, z)) \leftrightarrow (\text{Sat}(x, y) \wedge \text{Sat}(x, z)))$$

(S4)

$$\text{ass}(x, \text{all}(\ulcorner v \urcorner, y)) \rightarrow (\text{Sat}(x, \text{all}(\ulcorner v \urcorner, y)) \leftrightarrow (\forall z \supseteq x) \text{Sat}(z, y))$$

(S5)

$$\text{Sat}(x, y) \rightarrow \text{ass}(x, y)$$

We recall that CT is the theory obtained from CT^\uparrow from Def. 4.27. by allowing formulas of \mathcal{L}_T , in which the truth predicate only applies to formulas of \mathcal{L} , into the schema of Δ_0 -induction. As one might expect, st and tr give us

LEMMA 4.35.

- (i) CT defines the truth predicate of CS;
- (ii) CS defines the truth predicate of CT.

Proof. (i) crucially employs the properties of substitution. For instance, for S1, we reason as follows, with the contextual information that $\text{ass}(x, \ulcorner R(v_1, \dots, v_n) \urcorner)$ and $\bigwedge_{i=1}^n x(\ulcorner v_i \urcorner) = x_i$:

$$\begin{aligned} \text{Tsb}(x, \ulcorner R(v_1, \dots, v_n) \urcorner) &\leftrightarrow \text{T}^\ulcorner R(\dot{x}_1, \dots, \dot{x}_n) \urcorner && \text{by sb and EA} \\ &\leftrightarrow R(x_1, \dots, x_n) && \text{by Tat} \end{aligned}$$

For (S4), we notice that EA proves

$$(16) \quad \text{all}(\ulcorner v \urcorner, \text{sb}(x, y)) = \text{sb}(x, \text{all}(\ulcorner v \urcorner, y))$$

Now if $\text{Tsb}(x, \text{all}(\ulcorner v \urcorner, y))$, also $\text{Tall}(\ulcorner v \urcorner, \text{sb}(x, y))$ and therefore $\forall w \text{Tsub}(\text{sb}(x, y), \bar{w})$ by the CT axioms. If there is a mapping $z \supseteq x$ such that $\neg \text{Tsb}(z, y)$, also $\text{ass}(z, y)$ and $\text{ass}(x, \text{all}(\ulcorner v \urcorner, y))$. Therefore, there is a w_0 such that $\neg \text{dom}(x, w_0)$ and $\text{dom}(z, w_0)$. Hence, in EA,

$$(17) \quad \text{sub}(\text{sb}(x, y), \bar{w}_0) = \text{sb}(z, y)$$

We can then conclude $\text{Tsb}(z, y)$, contradicting our assumption. The converse direction is obtained similarly, with the help of (16) and (17).

(ii) is obtained once we have established a version of D above via suitable induction on the formal complexity of the ‘formula’ y :

$$(18) \quad \text{ass}(x, y) \rightarrow (\text{Sat}(x, y) \leftrightarrow \text{Sat}(\langle \rangle, \text{sb}(x, y)))$$

□

PROPOSITION 4.36. *CT and CS are truth equivalent.*

Proof. That CT is a retract of CS is immediate by definitions of st and tr. To conclude that CS is a retract of CT, we only need (18). In both cases the required isomorphisms are in fact identities. \square

We notice that, unlike the previous observations concerning disquotation, Proposition 4.36. still holds when we formulate CS in a language without domain constants. Moreover, since we essentially employ the extended induction of CS to obtain the t-equivalence of the two theories, it would be interesting to know what happens if we consider induction free versions of the theories involved. The question of whether $CS\upharpoonright$ and $CT\upharpoonright$ are t-equivalent (or mutually truth-definable) is in fact still open.

In desideratum (III) above, we required the notions of t-retracts and t-equivalence to be nonempty. Although not exciting, the examples just presented suffice to accomplish this minimal task. Admittedly, the main role of our proposed sufficient condition for conceptual equivalence is to discourage inadequate claims, and this is mostly accomplished by negative results such as Corollary 4.26..

I now move to an example of *bi-interpretability simpliciter* between truth theories that originates in (Leigh & Nicolai 2013), in which theories with external truth were first introduced.³⁰ Here I only sketch a simplified construction that may be useful to obtain more sophisticated – and more natural – examples. Let EA^2 be obtained by ‘cloning’ the language \mathcal{L} in a ‘two-sorted’ version: in practice, I work with relativizing predicates $s(x)$ and $t(x)$ that represent the two copies of our numbers. I abbreviate $\forall x(s(x) \rightarrow \varphi(x))$ with $(\forall x : s) \varphi(x)$, and similarly for t . The language of EA^2 will thus contain, for instance, two versions of the primitive predicate ‘... is zero’, Z_s and Z_t , and so on for the other primitive notions, including identity symbols $=_s, =_t$. The induction axioms of EA^2 come therefore in two flavours, one in which the relevant variable can only belong to s and one in which it can only belong to t , although formulas from the entire language are allowed in both types of instances. We assume the usual arithmetization of the syntax for EA^2 carried out in the s -portion of the language. We also include in the theory a ‘Frege relation’ F , that is a relation witnessing an isomorphism between the two domains corresponding to the two sorts. In other words we add axioms of the form:

- (F0) $x F y \rightarrow s(x) \wedge t(y)$
- (F1) ‘ F is a bijection between s -objects and t -objects’
- (F2) $(\forall x : s)(\forall y : t)(x F y \rightarrow (Z_s(x) \leftrightarrow Z_t(y)))$
- (F3) $(\forall x, u : s)(\forall y, v : t)(x F y \wedge u F v \rightarrow (S_s(x, u) \leftrightarrow S_t(y, v)))$

By formal induction in EA^2 , the analogues of F2 and F3 for the predicates A and M for addition, multiplication and exponentiation come out as theorems of EA^2 :

³⁰ Although already at that time the influential work of Richard Heck, who should be granted the priority, together with Albert Visser, of introducing external truth, was widely circulating. This approach can be motivated in the following way: If one starts with an arbitrary object theory U , unlike what it is commonly thought, it is not so straightforward to introduce a notion of truth for U . If, for instance, U fails to be sequential (cf. §1), the theory does not have a good notion of sequence and we won’t have the necessary resources to express full satisfaction in a direct or derivative form – e.g. by employing domain constants as in the case of a unary truth predicate. However, truth and sequences can be added ‘from the outside’ (Leigh & Nicolai 2013; Heck 2015; Nicolai 2016).

$$(19) \quad (\forall x, y, z : s)(\forall u, v, w : t)(xFu \wedge yFv \wedge zFw \rightarrow (A_s(x, y, z) \leftrightarrow A_t(u, v, w)))$$

$$(20) \quad (\forall x, y, z : s)(\forall u, v, w : t)(xFu \wedge yFv \wedge zFw \rightarrow (M_s(x, y, z) \leftrightarrow M_t(u, v, w)))$$

$$(21) \quad (\forall x, y : s)(\forall u, v : t)(xFu \wedge yFv \rightarrow (E_s(x, y) \leftrightarrow E_t(u, v)))$$

I then introduce *typed* truth axioms for EA^2 by employing a truth predicate Tr of type s , that is, applying only to objects of sort s . The characterizing truth axioms, besides the obvious analogues of T2 and (\star) (cf. p. 18), are the ones for atomic formulas and quantifiers:

$$(Tat^*) \quad Tr^\top R_s(\dot{x}_1, \dots, \dot{x}_n)^\top \leftrightarrow R_t(F(x_1), \dots, F(x_n)) \text{ for all } s\text{- and } t\text{-relation symbols}$$

$$(Tq) \quad Fml^1_{\mathcal{L}^s}(y) \rightarrow \left(Tr(\text{all}_s(\ulcorner v^\top, y \urcorner)) \leftrightarrow (\forall z : s) (Tr \text{ sub}_s(y, \bar{z})) \right)$$

I refer to the resulting theory of truth as $T[EA^2]$. The theory CT has been already introduced: it is convenient to think of the current version of CT as built over a one-sorted language of sort s , a sublanguage of the language of EA^2 .

PROPOSITION 4.37. $T[EA^2]$ and CT are bi-interpretable.

Proof.

I specify the interpretation $K: T[EA^2] \rightarrow CT$. The idea behind it is entirely straightforward: the two copies of the numbers in $T[EA^2]$ are reproduced in CT as pairs with 0 and 1 as first members, whereas primitive relations collapse into their counterparts in CT. In CT we find the following elementary formulas:

$$\text{pair}(x) :\leftrightarrow 'x \text{ is an ordered pair}'$$

$$\text{pri}(x) :\leftrightarrow 'x \text{ is an ordered pair with first member } i'$$

$$\pi_i(x, y) :\leftrightarrow 'y \text{ is the } i^{\text{th}} \text{ projection of the pair } x'$$

I employ $\pi_i(\cdot)$ in a functional form. Let

$$t^K(x) :\leftrightarrow x : \text{pr}0$$

$$s^K(x) :\leftrightarrow x : \text{pr}1$$

$$Z_t^K(x) :\leftrightarrow Z(\pi_1(x))$$

$$Z_s^K(x) :\leftrightarrow Z(\pi_1(x))$$

$$\vdots$$

$$\vdots$$

$$Tr^K_x :\leftrightarrow T\pi_1(x)$$

$$F^K(x, y) :\leftrightarrow x : \text{pr}1 \wedge y : \text{pr}0 \wedge \pi_1(x) = \pi_1(y)$$

$$x =_s y :\leftrightarrow x, y : \text{pr}1 \wedge \pi_1(x) = \pi_1(y) \quad x =_t y :\leftrightarrow x, y : \text{pr}0 \wedge \pi_1(x) = \pi_1(y)$$

K relativizes quantifiers to $\text{pair}(x)$. The vertical dots refer to the clauses for the other arithmetical relations (in both versions), and the convention employed in the definition of EA^2 concerning relativized quantifiers has been extended to $\text{pr}0$ and $\text{pr}1$.

The interpretation $L: CT \rightarrow T[EA^2]$ simply relativizes everything to the predicate s . Crucially:

$$(\forall x \phi)^L :\leftrightarrow (\forall x : s) \phi^L$$

$$T^L(x) :\leftrightarrow Tr(x)$$

By letting, in $T[EA^2]$,

$$G(x, y) :\leftrightarrow (\text{pr}1^s(x) \wedge \pi_1^s(x) = y) \vee (\text{pr}0^s(x) \wedge F(\pi_1^s(x)) = y),$$

it is not difficult to check that $G: L \circ K \cong \text{id}_{T[EA^2]}$ and $\pi_1: K \circ L \cong \text{id}_{CT}$, witnessing the required bi-interpretability of $T[EA^2]$ and CT .

□

We conclude this section by with a chart summarizing our findings.

| DESIDERATUM | EXAMPLES |
|---|---|
| Extensions of mutual truth-definability | t-retracts, t-equivalence |
| Separating t-retracts, t-equivalence, mutual truth-definability | PUTB, KF; PUTB _↓ , KF _↓ |
| Non-emptiness of t-retracts, t-equivalence | Flavours of Truth and Satisfaction |
| Bi-interpretability simpliciter | ‘External’ truth |

§5. Extending the framework: syntactical embeddings, e-retractions, e-equivalence

As anticipated in the initial section, natural generalizations of the notions of t-retract and t-equivalence may yield new insights on the comparison between the operations of adding typed truth and predicative comprehension to a ground syntactic structure. (Nicolai 2016) focuses on abstracting away restrictive choices of the object theory. I now consider a less general framework but stricter notions of reduction. In particular, claims like the following belong to the truth-theoretic and foundational folklore:

Typing truth predicates corresponds to a much more severe move in the case of comprehension: typing corresponds to predicative typed comprehension [...] Actually ramified type theory over Peano arithmetic as base theory, which is known as ramified analysis, is *equivalent* to typed compositional truth. (p. 28)

In what sense should this ‘equivalence’ be understood? Most likely Halbach refers to *proof-theoretical equivalence* (cf. (Feferman 1988)), which in the cases I consider below is no different from a more liberal notion of truth-definition. In this section I employ analogues of the notions of truth-definitions, t-retractions and t-equivalence to refine these folklore statements. Several observations contained in this section rely on unpublished work of Ali Enayat and Albert Visser.

I first generalize the notion of a truth-definition to arbitrary nonlogical vocabulary extending the language of a suitable syntax theory B .

DEFINITION 5.38. (SYNTACTICAL EMBEDDING) *Let $\mathcal{L}_B \subseteq \mathcal{L}_T, \mathcal{L}_W$ and $B \subseteq T, W$. Then we say that T is syntactically embeddable in W if and only if there is a relative interpre-*

tation $K: T \rightarrow W$ that leaves the primitive vocabulary of \mathcal{L}_B unchanged and does not (non-trivially) relativize its quantifiers.

Truth-definitions as defined above are clearly examples of syntactical embeddings. We now move to generalizations of the notions of t-retract and t-equivalence.

DEFINITION 5.39. (E-RETRACT, E-EQUIVALENCE) *Let $T, W \supseteq B$ be given.*

- (i) *T is an e-retract of W if there are syntactical embeddings $K: T \rightarrow W$ and $L: W \rightarrow T$ and a T -definable $F: L \circ K \cong \text{id}_T$;*
- (ii) *T and W are e-equivalent if there are syntactical embeddings $K: T \rightarrow W$ and $L: W \rightarrow T$, a T -definable $F: L \circ K \cong \text{id}_T$ and a W -definable $G: K \circ L \cong \text{id}_W$.*

To be able to apply the new notions, I briefly recall some standard definitions concerning subsystems of second-order arithmetic as they can be found, for instance, in (Simpson 2009). \mathcal{L}_2 is the two-sorted language extending \mathcal{L} with a sort for sets of natural numbers, or ‘reals’. The two kinds of variables will often be denoted by x, y, z, \dots and X, Y, Z, \dots respectively, possibly with indices. In practice, to conform with our notions of reduction, it is convenient to consider the two-sorted language as notational abbreviation for a language with suitable relativizing predicates. That is, we officially work in a single-sorted language with relativizing predicates $\text{se}(x)$, ‘ x is a set’, and $\text{nu}(x)$, ‘ x is a number’, and a primitive membership predicate \in . A formula of \mathcal{L}_2 is said to be *arithmetical* if it contains no set quantifiers but possibly set parameters. The so-called *arithmetical comprehension schema* has the form

$$(\text{aca}) \quad \exists X \forall x (x \in X \leftrightarrow \varphi(x, \vec{u}, \vec{Y}))$$

where $\varphi(v_0)$ is arithmetical and does not contain X .³¹ The theory ACA in \mathcal{L}_2 is obtained by relativizing the basic axioms of EA to the ‘numbers’ sort and by extending it with (aca) and an induction schema for arbitrary \mathcal{L}_2 -formulas. ACA^{pf} in \mathcal{L}_2 is obtained from ACA by disallowing set-parameters into instances of the comprehension schema. The theory ACA_0 is obtained from ACA by replacing its full induction schema with the equivalent, in the official language, of the single sentence

$$\text{IND} \quad (0 \in X \wedge \forall x (x \in X \rightarrow Sx \in X)) \rightarrow \forall x x \in X$$

In the official presentation of the theories by means of relativizing predicates, I also require the domain to be partitioned by them.

For the purpose of this section I will consider only *finite* iterations of full predicative comprehension obtained intuitively by iterating ACA n -times. The restriction to finite levels is motivated by the difficulties in the formulation of iterations of predicative comprehension for higher ordinals. In particular, ramified analysis is classically formulated by employing reflection principles (see (Feferman 1964)), and this creates some difficulties in transferring the results below to the transfinite.

More formally, we have

$$\begin{aligned} \mathcal{L}^{<0} &:= \mathcal{L} \\ \mathcal{L}^{<n+1} &:= \mathcal{L} \cup \{\text{se}^i, \in^i\} \end{aligned} \quad \text{for } i \leq n$$

³¹ In our official formulation, the arithmetical comprehension thus becomes

$$\text{se}(\vec{z}) \rightarrow \exists y (\text{se}(y) \wedge \forall u (\text{nu}(u) \rightarrow (u \in y \leftrightarrow \varphi(u, \vec{v}, \vec{z}))))$$

I will also write \mathcal{L}^n for $\mathcal{L}^{<n+1}$. RA_0 is ACA itself. The *degree* of a formula of \mathcal{L}^n is the maximum of the $k+1$ such that $\text{se}^k(y)$ appears in the formula with y bound and of the m such that $\text{se}^m(z)$ appears in the formula with z is free. I write φ^n for a formula of at most degree n . Then RA_{n+1} results from RA_n by extending its induction schema to the new language and by adding to it the axioms:

$$(22) \quad x \in^{n+1} y \rightarrow \text{se}^{n+1}(y)$$

$$(23) \quad \exists y (\text{se}^{n+1}(y) \wedge \forall x (x \in^{n+1} y \leftrightarrow \varphi^{n+1}(x, \vec{u}, \vec{X}^i)))$$

For reasons of readability, we have omitted the relativization $\text{nu}(x)$. We will keep this convention in what follows. In (23) y is not free in φ and \vec{X}^i is a string of parameters from elements of se^i for $i \leq n+1$. By disallowing (set-)parameters into the schema (23), one obtains a parallel hierarchy of finitely iterated (parameter-free) predicative comprehension, with $\text{RA}_0^{\text{pf}} := \text{ACA}^{\text{pf}}$.

We are interested in associating the theories just presented with finite iterations of $\text{UTB}[\text{PA}]$ and $\text{CT}[\text{PA}]$. The latter theories are formulated in the languages $\mathcal{L}_T^n := \mathcal{L} \cup \{T_0, \dots, T_n\}$ for $n \in \omega$, with $\mathcal{L}_T^{<0} := \mathcal{L}$ and $\mathcal{L}_T^{<n+1} := \mathcal{L}_T^n$.

DEFINITION 5.40. Let $\text{RT}_0 := \text{CT}[\text{PA}]$ and $\text{RDT}_0 := \text{UTB}[\text{PA}]$, both formulated in \mathcal{L}_T^0 .

- The theory RDT_{n+1} is obtained by extending RDT_n with full \mathcal{L}_T^{n+1} -induction and all instances of

$$\forall x (T_{n+1} \ulcorner \varphi(\dot{x}) \urcorner \leftrightarrow \varphi(x))$$

for all \mathcal{L}_n -formulas $\varphi(v)$.

- The theory RT_{n+1} contains the axioms of RT_n , including the induction schema for the entire language \mathcal{L}_T^{n+1} , and the following, for $m < n+1$:

$$(R1_{n+1}) \quad \forall x_1, \dots, x_n (T_{n+1} \ulcorner R(\dot{x}_1, \dots, \dot{x}_n) \urcorner \leftrightarrow R(x_1, \dots, x_n))$$

$$(R2_{n+1}) \quad \forall x (\text{Sent}_{\mathcal{L}_T^{<m}}(x) \rightarrow (T_{n+1} \ulcorner T_m \dot{x} \urcorner \leftrightarrow T_m x))$$

$$(R3_{n+1}) \quad \forall x (\text{Sent}_{\mathcal{L}_T^n}(x) \rightarrow (T_{n+1} \text{ng}(x) \leftrightarrow \neg T_{n+1} x))$$

$$(R4_{n+1}) \quad \forall x, y (\text{Sent}_{\mathcal{L}_T^n}(\text{and}(x, y)) \rightarrow (T_{n+1} \text{and}(x, y) \leftrightarrow (T_{n+1} x \wedge T_{n+1} y)))$$

$$(R5_{n+1}) \quad \forall x, v (\text{Sent}_{\mathcal{L}_T^n}(\text{all}(v, x)) \rightarrow (T_{n+1} \text{all}(v, x) \leftrightarrow \forall y T_{n+1} \text{sub}(x, \bar{y})))$$

$$(R6_{n+1}) \quad \forall y \leq n \forall x (\text{Sent}_{\mathcal{L}_T^{<y}}(x) \rightarrow (T_{n+1} \ulcorner T_y \dot{x} \urcorner \leftrightarrow T_{n+1} x))$$

The standard reductions between ramified truth and predicative comprehension mainly amount to mutual syntactical embeddings (cf. (Halbach 2014; Feferman 1991; Takeuti 1987; Fischer 2009)): Let $T_0(X^0, u)$ be the \mathcal{L}_0 -formula

$$\begin{aligned} & \forall y (y \in X^0 \leftrightarrow \text{Sent}_{\mathcal{L}}(y) \wedge \text{lc}(x) \leq u \wedge \\ & \quad \forall \vec{x} (y = \ulcorner R(\dot{x}) \urcorner \rightarrow (y \in X^0 \leftrightarrow R\vec{x}))) \\ & \quad \forall z (y = \text{ng}(z) \rightarrow (y \in X^0 \leftrightarrow z \notin X^0)) \wedge \\ & \quad \forall w, z (y = \text{and}(w, z) \rightarrow (y \in X^0 \leftrightarrow (w \in X^0 \wedge z \in X^0))) \wedge \\ & \quad \forall v, z, y (x = \text{all}(v, z) \rightarrow (x \in X^0 \leftrightarrow \forall y (\text{sub}(z, \bar{y}) \in X^0))) \end{aligned}$$

Moreover, let $\mathbb{V}_0(x) := (\exists Y^0)(T_0(Y^0, \text{lc}(x)) \wedge x \in Y^0)$. This latter formula plays the role of a partial truth predicate for \mathcal{L} -formulas of complexity up to (and including) the complexity

of x . To define what it means to be a partial truth set of level $n+1$, we generalize the notion of logical complexity of a formula to sentences of \mathcal{L}_T^n by letting $\text{lc}^{n+1}(\ulcorner T_n \dot{y} \urcorner)$ to be 0 and keeping the rest as it is:

$$\begin{aligned} \mathbb{T}_{n+1}(X^{n+1}, u) :& \leftrightarrow \forall y (y \in X^{n+1} \leftrightarrow \text{Sent}_{\mathcal{L}_T^n}(y) \wedge \text{lc}^{n+1}(y) \leq u \wedge \\ & \forall \vec{x} (y = \ulcorner R(\vec{x}) \urcorner \rightarrow (y \in X^{n+1} \leftrightarrow R\vec{x})) \wedge \\ & \forall z (y = \text{ng}(z) \rightarrow (y \in X^{n+1} \leftrightarrow z \notin X^{n+1})) \wedge \\ & \forall w, z (y = \text{and}(w, z) \rightarrow (y \in X^{n+1} \leftrightarrow (w \in X^{n+1} \wedge z \in X^{n+1}))) \wedge \\ & \forall v, z, y (x = \text{all}(v, z) \rightarrow (x \in X^{n+1} \leftrightarrow \forall y (\text{sub}(z, \bar{y}) \in X^{n+1}))) \wedge \\ & \forall z (\text{Sent}_{\mathcal{L}_T^{<n}}(z) \wedge y = \ulcorner T_n \dot{z} \urcorner \rightarrow (y \in X^{n+1} \leftrightarrow \mathbb{V}_n(z))) \wedge \\ & \forall u \leq n \forall z (\text{Sent}_{\mathcal{L}_T^{<u}}(z) \wedge y = \ulcorner T_u \dot{z} \urcorner \rightarrow (y \in X^{n+1} \leftrightarrow z \in X^{n+1}))) \end{aligned}$$

The formula $\mathbb{V}_{n+1}(x)$ is then defined in the obvious way.

We finally adapt the folklore translations to our setting by defining the translation functions $K_n: \mathcal{L}_T^n \rightarrow \mathcal{L}^n$, $L_n: \mathcal{L}^n \rightarrow \mathcal{L}_T^n$ in Table 1. We remark that \mathcal{L}_n and \mathcal{L}_T^n feature only one identity symbol.

| | | |
|---|---|-------------------------|
| $R^{K_n}(\vec{x}) : \leftrightarrow R\vec{x}$ | $R^{L_n}(\vec{x}) : \leftrightarrow R\vec{x}$ | for $R \in \mathcal{L}$ |
| $x =^{K_n} y : \leftrightarrow x = y$ | $x =^{L_n} y : \leftrightarrow x = y$ | |
| $T_m^{K_n} x : \leftrightarrow \mathbb{V}_m(x)$ | $x(\in^m)^{L_n} y : \leftrightarrow T_m \text{sub}^1(x, y)$ | with $m \leq n$ |
| $(\neg \varphi)^{K_n} : \leftrightarrow \neg \varphi^{K_n}$ | $(\neg \varphi)^{L_n} : \leftrightarrow \neg \varphi^{L_n}$ | |
| $(\varphi \wedge \psi)^{K_n} : \leftrightarrow \varphi^{K_n} \wedge \psi^{K_n}$ | $(\varphi \wedge \psi)^{L_n} : \leftrightarrow \varphi^{L_n} \wedge \psi^{L_n}$ | |
| $(\forall x \varphi)^{K_n} : \leftrightarrow (\forall x : \text{nu}) \varphi^{K_n}$ | $(\forall x \varphi)^{L_n} : \leftrightarrow \forall x \varphi^{L_n}$ | |
| | $\text{nu}^{L_n}(x) : \leftrightarrow x = x$ | |
| | $(\text{se}^m)^{L_n}(x) : \leftrightarrow \text{Fml}_{\mathcal{L}^{<m}}^1(x)$ | $m \leq n$ |

Table 1. *The translations K_n, L_n .*

In the table, the elementary formula $\text{sub}^1(x, y) = z$ expresses the result of formally replacing the single free variable of a formula y with x . The folklore reductions can be summarized as follows.

PROPOSITION 5.41.

- (i) For each $n \in \omega$, RDT_n and RA_n^{pf} are mutually syntactically embeddable.
- (ii) For each $n \in \omega$, RT_n and RA_n are mutually syntactically embeddable.

Proof. One can check that K_n and L_n are the required syntactical embeddings, given the following facts:

- (24) $RA_n^{\text{pf}} \vdash \exists X^m \mathbb{T}_m(X^m, \bar{k})$ for $m \leq n$ and $k \in \omega$
- (25) $RA_n \vdash \forall x \exists X^n \mathbb{T}_n(X^n, x)$
- (26) if $\varphi(v)$ is in \mathcal{L}^m , then $\varphi^{L^m}(v)$ belongs to \mathcal{L}_T^m
- (27) $RT_n \vdash \forall u, \vec{z} (\mathbb{T}_m \ulcorner \varphi^{L^n}(\dot{u}, \vec{z}) \urcorner \leftrightarrow \varphi^{L^n}(u, \vec{z}))$ for $m \leq n$, φ in $\mathcal{L}^{<m}$, and $\text{se}^{<m}(\vec{z})$

□

From the previous chapter we know that e-retractions and e-equivalences properly strengthen mutual syntactical embeddings. It is thus interesting to investigate whether the folklore mutual embeddings between typed truth and comprehension can be lifted to stricter notions of equivalence. The following shows that the folklore translations partially suffice.

PROPOSITION 5.42. *For each $n \in \omega$, RT_n is an e-retract of RA_n .*

Proof. The case $n = 0$ follows from the fact that in $CT[PA]$ we have partial truth predicates $\text{Tr}_{\Sigma_n}(\cdot)$ of complexity Σ_n for Σ_n -formulas and for each n . I.(d). By its extended induction, we can prove in $CT[PA]$ that

$$(28) \quad T^{L_0 \circ K_0}(x) \leftrightarrow \exists y T_0 \ulcorner \text{Tr}_{\Sigma_0}(\dot{x}) \urcorner \leftrightarrow T_0 x$$

We emphasize that in (28) the truth predicate T_0 enables us to quantify over the index of the partial truth predicate. The required isomorphism is thus the identity on RT_0 .

The same idea can be extended to the languages \mathcal{L}_n , by focusing now on the hierarchy of Σ_n^T -formulas, and constructing in the standard way partial truth predicates $\text{Tr}_{\Sigma_n^T}(x)$ for the languages \mathcal{L}_T^n ; they should be thought as formalizing the stages of the construction of T_{n+1} from previously defined truth predicates. In particular, all \mathcal{L}_T^n -sentences φ deemed true at previous stages are such that $T_{n+1} \ulcorner \text{Tr}_{\Sigma_0^n}(\ulcorner \varphi \urcorner) \urcorner$ is provable in RT_{n+1} .

Again by full \mathcal{L}_T^{n+1} -induction:

$$(29) \quad RT_{n+1} \vdash T^{L_{n+1} \circ K_{n+1}}(x) \leftrightarrow \exists y T_{n+1} \ulcorner \text{Tr}_{\Sigma_n^T}(\dot{x}) \urcorner \leftrightarrow T_{n+1}(x)$$

As before, the required isomorphism is indeed the identity function on RT_{n+1} . □

But Proposition 5.42. is in a sense the best we can do. I adapt to the present setting an unpublished argument by Enayat and Visser and show that ramified analysis, both in full or parameter-free form, cannot be a retract of ramified truth, both in full compositional or uniform disquotational form.³²

PROPOSITION 5.43.

- (i) RA_n^{pf} is not a retract nor an e-retract of $RDT_n[PA]$, for $n \in \omega$.

³² Note that the argument is so general that can be applied also to the transfinite, once a suitable formulation of ramified analysis is given. I thank Albert Visser for giving me the permission of quoting his unpublished work.

(ii) RA_n is not a retract nor an e -retract of $RT_n[PA]$.

Proof. I give the argument for (i); the same reasoning applies to (ii). If RA_n^{pf} were a retract of RDT_n , any model $(\mathcal{M}_0, \mathcal{R}_0, \dots, \mathcal{R}_n)$ of RA_n^{pf} – where the \mathcal{R}_i 's are subsets of \mathcal{M}_0 – could define an internal model $(\mathcal{M}_1, \mathcal{S}_0, \dots, \mathcal{S}_n) \models RDT_n$ of \mathcal{M}_0 that would in turn define within itself an internal model $(\mathcal{M}_2, \mathcal{T}_0, \dots, \mathcal{T}_n) \models RA_n^{pf}$ of \mathcal{M}_1 such that \mathcal{M}_0 is isomorphic to \mathcal{M}_2 definably in \mathcal{M}_0 (cf. Figure 1, §3).

Let us now start with the ‘standard’ model of RA_n^{pf} , i.e. the tuple $(\omega, \mathcal{P}(\omega))$. By the above considerations, we have the situation depicted in Figure 2, where $(\mathcal{M}_1, \vec{\mathcal{T}}) \models RDT_n$ and $(\mathcal{M}_2, \vec{\mathcal{T}}) \models RA_n^{pf}$ – obviously the \mathcal{S}_i are subsets of $\text{Sent}_{\mathcal{L}^{<n}}$ and the \mathcal{T}_i 's subsets of \mathcal{M}_2 . Since $(\mathcal{M}_2, \vec{\mathcal{T}}) \cong (\omega, \mathcal{P}(\omega))$, there is an interpretation of $(\omega, \mathcal{P}(\omega))$ in $(\mathcal{M}_1, \vec{\mathcal{T}})$.

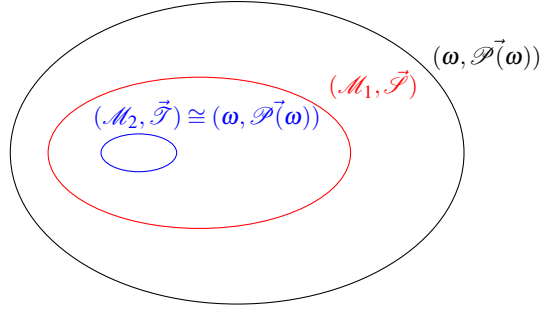


Fig. 2

Thus $(\mathcal{M}_1, \vec{\mathcal{T}})$ defines its standard natural numbers. At the same time, since $(\mathcal{M}_1, \vec{\mathcal{T}})$ satisfies induction with the truth predicate and it interprets the basic axioms of EA, by Lemma 2.2. one can $(\mathcal{M}_1, \vec{\mathcal{T}})$ -define an injection $f: \mathcal{M}_1 \rightarrow \omega$ by primitive recursion. Therefore \mathcal{M}_1 is countable, and thus it cannot define the uncountable structure $(\omega, \mathcal{P}(\omega))$. \square

COROLLARY 5.44.

- (i) For each $n \in \omega$, RA_n^{pf} and RDT_n are not e -equivalent nor bi-interpretable;
- (ii) For each $n \in \omega$, RA_n and RT_n are not e -equivalent nor bi-interpretable;

COROLLARY 5.45. *The relations ‘being an e -retract of’ and e -equivalence are properly stricter than mutual syntactical embeddability.*

Corollary 5.44. is already worth consideration. It denies the possibility of strong forms of equivalence between predicative comprehension and typed truth. Typed compositional truth predicates seem to be ‘stronger’, in the precise sense described above, than the corresponding set-theoretic axioms. This gives us already a refinement of the folklore claim of the correspondence between typed truth and predicative comprehension, at least in the case of full compositional truth. A natural question therefore is to investigate the nature of this asymmetry by measuring what exactly one has to add to ramified analysis to recapture the corresponding truth predicates in the sense of e -equivalence. To do this, at least in the case of RA_n and RT_n , I again elaborate on an idea by Enayat and Visser.

I move to definitional extensions of the RA_n 's obtained by adding identity predicates $=_m$, for $m \leq n$, satisfying

$$(30) \quad y =_m z \leftrightarrow \forall u (u \in^m y \leftrightarrow u \in^m z)$$

$$(31) \quad y =_m z \rightarrow \text{se}^m(y) \wedge \text{se}^m(z)$$

For simplicity, I again employ RA_n and RA_n^{pf} as names of the extended theories. Also I keep K_m from Table 1 fixed. In RT_n , for $m \leq n$, I define

$$(32) \quad x \sim_m y :\leftrightarrow \text{Fml}_{\mathcal{L}_T^{<m}}^1(x) \wedge \text{Fml}_{\mathcal{L}_T^{<m}}^1(y) \wedge \forall u (\text{T}_m \text{sub}^1(u, x) \leftrightarrow \text{T}_m \text{sub}^1(u, y))$$

(32) simply says that the $\mathcal{L}_T^{<m}$ -formulas x, y are satisfied by the same objects. The translation M_m is obtained from L_m in Table 1 by introducing clauses for $=_m$ as follows, with $k \leq m$:

$$(33) \quad x =_k^M y :\leftrightarrow x \sim_k y$$

M_n determines a syntactical embedding of RA_n and RA_n^{pf} in RA_n and RDT_n respectively. Moreover, the argument in Proposition 5.42. proceeds unchanged if M_n is employed instead of L_n . Let

$$D_m :\leftrightarrow \bigwedge_{i=0}^m \forall x \left(\text{se}^i(x) \rightarrow \exists y \left(\text{Fml}_{\mathcal{L}_T^{<i}}(y) \wedge \forall u (u \in^i x \leftrightarrow \forall_i(\text{sub}^1(u, y))) \right) \right)$$

The sentence D_n says that for all $i \leq n$, every set of level i can be defined by a $\mathcal{L}_T^{<i}$ -formula: by 'define' here I mean that its elements can be shown to be exactly the elements satisfying this formula. By reflecting on the proof of Proposition 5.43., one can immediately see how the addition of D_n excludes what we called the 'standard model of RA_n ' from the space of the models of $RA_n + D_n$. Moreover, D_n forces the $\mathcal{L}_T^{<i}$ -formula defining a set X^i to be unique up to $\sim_i^{K_i}$ -equivalence. By employing essentially the same reasoning as in Proposition 5.42., we have

LEMMA 5.46. *For each $n \in \omega$,*

- (i) M_n is a syntactical embedding of $RA_n + D_n$ in RT_n ;
- (ii) RT_n is an e-retract of $RA_n + D_n$.

Finally we can check that the sentences D_n suffice to restore the symmetry between the two hierarchies:

PROPOSITION 5.47. *$RA_n + D_n$ is an e-retract of RA_n .*

Proof. We know that $RA_n + D_n$ and RT_n are mutually syntactically embeddable via K_n and M_n . The induction schema of $RA_n + D_n$ enables one to prove, for all $m \leq n$ (and for all $n \in \omega$),

$$(34) \quad \forall y \left(\text{Fml}_{\mathcal{L}_T^{<m}}(y) \rightarrow (\exists X^m)(\forall u)(u \in^m X^m \leftrightarrow \forall_m(\text{sb}(u, y))) \right)$$

Next I let

$$(35) \quad H_n(x, y) :\leftrightarrow \bigvee_{i=0}^n (\text{Fml}_{\mathcal{L}_T^{<i}}(x) \wedge \text{se}^i(y) \wedge (\forall u)(u \in^i y \leftrightarrow \forall_i(\text{sub}^1(u, x)))) \vee (\neg \text{Fml}_{\mathcal{L}_T^{<n}}^1(x) \wedge y = x)$$

It remains to verify, in $RA_n + D_n$, that $H_n(x, y)$ is the required isomorphism from $K_m \circ M_m$ to $\text{id}_{RA_n + D_n}$, that is, that conditions (2)-(9) on page 10 are satisfied by H . D_n , in combination with (34), give us the totality conditions for H . I now verify that $=_i$ and $(\in^i)^{K_n \circ M_n}$ behave as expected for $i \leq n$. This will complete the proof.

We first notice that $H_n(x, y)$ is ‘functional’, that is, for $m \leq n$:

$$(36) \quad H_m(x, y) \wedge H_m(x, z) \rightarrow y =_m^{K_m \circ M_m} z$$

$$(37) \quad H_m(x, y) \wedge H_m(z, y) \rightarrow x =_m z$$

Assuming $H_i(u, w)$ and $H_i(x, y)$ with $\text{Fml}_{\mathcal{L}_T^{<i}}(x)$, I first show that

$$RA_n + D_n \vdash H_i(u, w) \wedge H_i(x, y) \rightarrow (u =_i^{K_n \circ M_n} x \leftrightarrow w =_i y)$$

We can safely assume that x, u are $\mathcal{L}_T^{<i}$ -formulas. We have

$$\begin{aligned} u =_i^{K_n \circ M_n} x &\leftrightarrow u \sim_i^{K_n} x \\ &\leftrightarrow \forall v (\nabla_i(\text{sub}^1(v, u)) \leftrightarrow \nabla_i(\text{sub}^1(v, x))) \\ &\leftrightarrow \forall v (v \in^i w \leftrightarrow v \in^i y) \\ &\leftrightarrow w =_i y \end{aligned}$$

The penultimate line is obtained by definition of H_i . For \in^i I show

$$RA_n + D_n \vdash H_i(u, w) \wedge H_i(x, y) \rightarrow (u(\in^i)^{K_n \circ M_n} x \leftrightarrow w \in^i y)$$

Assuming that $H_i(u, w)$ and $H_i(x, y)$ we have, with $\text{Fml}_{\mathcal{L}_T^{<i}}(x)$ and $\text{nu}(u)$:

$$\begin{aligned} u(\in^i)^{K_n \circ M_n} x &\leftrightarrow T_i^{K_n}(\text{sb}(u, x)) \\ &\leftrightarrow \nabla_i(\text{sb}(u, x)) \\ &\leftrightarrow u \in^i y \\ &\leftrightarrow w \in^i y \end{aligned}$$

In the last line we have employed the fact that $\neg \text{Fml}_{\mathcal{L}_T^{<i}}(u)$. □

COROLLARY 5.48. $RA_n + D_n$ and RT_n are e -equivalent.

Finding principles that render iterated uniform disquotation equivalent to ramifications of parameter free comprehension is, as one might expect, more difficult. At least if one wants resort to principles that represent meaningful restrictions such as the D_n ’s above. For instance, already to achieve an analogue of Proposition 5.42. by using the same translations, we would require RDT_n to be able to prove a principle like (29). Although RDT_n has full induction, however, there seems to be no way to mimic the crucial role that compositional axioms have in its proof in RT_n . Of course we can add what we lack to RDT_n .

Let

$$E_n := \bigwedge_{i=0}^n \forall x (T_i x \leftrightarrow (\exists y : \text{Fml}_{\mathcal{L}_T^{<i}})(\mathbb{T}_i^{M_n}(y, \text{lc}^i(x)) \wedge T_i \text{sb}^1(y, x)))$$

It is somewhat tedious, although straightforward by E_n , to verify that, for $n \in \omega$,

LEMMA 5.49. $RDT_n + E_n \vdash E^{M_n \circ K_n}$.

Therefore $\text{RDT}_n + E_n$ and $\text{RA}_n^{\text{pf}} + E_n^{\text{Kn}}$ are mutually syntactically embeddable. With Lemma 5.49. at hand, we can conclude:

PROPOSITION 5.50. *$\text{RDT}_n + E_n$ is a retract of $\text{RA}_n^{\text{pf}} + E_n^{\text{Kn}}$.*

By adding new principles to the theories in Proposition 5.50., we might even be able to prove the e-equivalence of extensions of $\text{RDT}_n + E_n$ and $\text{RA}_n^{\text{pf}} + E_n^{\text{Kn}}$, although the more we add, the more the theories will look convoluted and hardly justifiable.

§6. Conclusion I have argued that mutual truth-definability, let alone mutual interpretability, is not adequate as sufficient condition for two theories of truth to embody the same conception of truth as defined in §1 and §2. The most prominent example of this failure is represented by case of the well-known principles of truth corresponding to the theories KF and PUTB, that are mutually truth-definable and still do not share most of the distinctive features described in §2. An alternative is the notion of t-equivalence: what I called Thesis 1 claims that it is a sufficient condition for the conceptual equivalence of the notions of truth arising from t-equivalent theories. To support the plausibility of this thesis, already suggested by the relationships between t-equivalence and synonymy given by Friedman and Visser's theorem, I have shown that it captures the manifest non-equivalence of the notions of truth of KF and PUTB – as the theories turn out to be not t-equivalent – and that some intuitively very close notions of truth, such as variations of Tarskian truth and satisfaction in the presence of full induction, result in fact in t-equivalent theories. T-retracts and t-equivalence, therefore, have been devised to be in continuity and extend mutual truth-definability, but if one loosens the criteria on the interpretation of the base-theoretic language, interesting combinations may arise. I have also provided in §4.2. a simple-minded template for generating theories of truth that are bi-interpretable but cannot be t-equivalent.

Mutual truth-definability is a particular case of a mutual syntactical embedding. By following the same strategy adopted to define t-retracts and t-equivalence, one can define the more general notions of e-retract and e-equivalence. These notions have proved to be useful to refine the folklore claims that link iterations of ramified truth and iterations of arithmetical comprehension. By crucially employing recent insights due Albert Visser and Ali Enayat, I have verified that finite iterations of ramified truth are not e-equivalent to finite iterations of arithmetical comprehension; this also shows that the notions of e-retract and e-equivalence are properly stricter than mutual syntactical embeddability. Since typed (disquotational or compositional) truth is an e-retract of finite iterations of arithmetical comprehension, the philosopher interested in deflating the ontological assumptions on the existence of subsets of \mathbb{N} may welcome this phenomenon as a confirmation of this possibility. At the same time, the failure of the t-equivalence between the two hierarchies should suggest that, from the logical point of view, typed truth and membership to (a portion of the) predicatively definable sets are *not* notational variants of each other.

Obviously these are only small, initial steps into the application to theories of truth and subsystems of analysis of the categories of theories and interpretations studied in (Visser 2006) and there are a number of unresolved issues. We list a few. In the first place one may consider more examples of natural theories of truth that are mutually truth-definable and investigate their t-equivalence; a natural starting point is, for instance, to consider the hierarchy RT_n for $n \in \omega$ and the corresponding (in the sense of mutual truth-definability) ω -consistent subsets of FS. A second example, related to bi-interpretability simpliciter, concerns theories of truth with the same truth theoretic axioms but built on bi-interpretable

base theories: the question is then whether they are uniformly bi-interpretable both in the typed and in the type-free case. In comparing truth axioms and arithmetical comprehension, even more variations are possible: one question is to consider theories with restricted induction (both in the truth-theoretic and the second-order side) that escape Proposition 5.43. and investigate whether they are e-retracts of each other and/or whether they are e-equivalent. Moreover, it is natural to wonder whether it is possible to lift the observations in §5. to transfinite iterations.

Acknowledgments I thank Ali Enayat, Kentaro Fujimoto, Volker Halbach, Graham Leigh, Lavinia Picollo, Johannes Stern, and the audiences in workshops in Oxford, Paris, Florence, for their feedback. Special thanks go to Albert Visser for allowing me to cite his unpublished arguments, to Thomas Schindler for correcting some inaccuracies in the definitions of the ramified truth theories in §5, and to the anonymous referees for their detailed reports.

BIBLIOGRAPHY

- Cantini, A. (1989). Notes on Formal Theories of Truth. *Zeitschrift für Mathematische Logik und Grundlagen der Mathematik* 35: 97–130.
- Cantini, A. (1990). A theory of truth formally equivalent to ID_1 . *The Journal of Symbolic Logic* 55.
- Cieśliński, C. (2010). Truth, Conservativeness, and Provability. *Mind* 474: 409–422.
- In Achourioti, D., J.M. Fernández, H. Galinon and K. Fujimoto (eds.). *Unifying the Philosophy of Truth*. Springer.
- Davidson, D. (1984). *Inquiries into Truth and Interpretation* Oxford University Press, Oxford.
- Enayat, A. and A. Visser (2015). New Constructions of Satisfaction Classes. In Achourioti, D., J.M. Fernández, H. Galinon and K. Fujimoto (eds.). *Unifying the Philosophy of Truth*. Springer.
- Enayat A., J. H. Schmerl, and A. Visser (2010). ω -models of finite set theory. In *Set theory, arithmetic, and foundations of mathematics: theorems, philosophies*, Juliette Kennedy, and Roman Kossak, eds. Lect. Notes Log., vol. 36, Assoc. Symbol. Logic, La Jolla, CA, 2010, pp. 43–65.
- Feferman, S. (1960). Arithmetization of Metamathematics in a General Setting. *Fundamenta Mathematicae*.
- Feferman, S. (1964). Systems of Predicative Analysis. *The Journal of Symbolic Logic* 29(1): 1–30.
- Feferman, S. (1991). Reflecting on Incompleteness. *The Journal of Symbolic Logic* 56: 1–49.
- Feferman, S. (1998). What rests on what? The proof-theoretic analysis of mathematics. In *In the Light of Logic*, Oxford University Press, Oxford, 187–208.
- Field, H. (2008). *Saving Truth from Paradox*. Oxford University Press.
- Fischer, M. (2009). Minimal Truth and Interpretability. *The Review of Symbolic Logic* 2(4) 2009.
- Fischer, M., V. Halbach, J. Kriener and J. Stern (2015). Axiomatizing semantic theories of truth?. *The Review of Symbolic Logic*.
- Friedman, H. and M. Sheard (1987). An axiomatic approach to self-referential truth. *Annals of Pure and Applied Logic* 33: 1–21.
- Friedman, H. and A. Visser (2014). When Bi-Interpretability Implies Synonymy. *Logic Group Preprint Series*, University of Utrecht.

- Fujimoto, K. (2010). Relative Truth Definability of Axiomatic Truth Theories. *The Bulletin of Symbolic Logic* 16(3): 305–344.
- Hájek P. and P. Pudlák (1998). *Metamathematics of First-Order Arithmetic*. Springer.
- Halbach, V. (1999). Disquotationalism and Infinite Conjunctions. *Mind* 108, 1–22.
- Halbach, V. (2009). Reducing compositional to disquotational truth. *The Review of Symbolic Logic* 2, 786–798.
- Halbach, L. (2014). *Axiomatic Theories of Truth*. Revised Edition. Cambridge University Press.
- Halbach, V. and L. Horsten (2006). Axiomatizing Kripke’s Theory of Truth. *The Journal of Symbolic Logic* 71: 667–712.
- Halbach, V. and L. Horsten (2015). Norms for Theories of Reflexive Truth. In Achourioti, D., J.M. Fernández, H. Galinon and K. Fujimoto (eds.). *Unifying the Philosophy of Truth*. Springer.
- Halbach, V. and C. Nicolai (2016). *On the Costs of Nonclassical Logic*. Manuscript.
- Heck, R. (2015) Consistency and the Theory of Truth. *The Review of Symbolic Logic*
- Horsten, L. (2011). *The Tarskian Turn: deflationism and axiomatic truth*. MIT Press, Cambridge MA.
- Kaye, R. (1991). *Models of Peano Arithmetic*. Oxford University Press.
- Kaye, R. & T. L. Wong (2007) On interpretations of arithmetic and set theory. *Notre Dame Journal of Formal Logic* 48(4): 497–510, 2007.
- Kripke, S. (1975). Outline of a theory of truth. *Journal of Philosophy* 72: 690–712.
- Laurence, S. and E. Margolis (1999). *Concepts: Core Readings*. Cambridge, MA: MIT Press.
- Leigh, G. E. (2015). Conservativity for theories of compositional truth via cut elimination, to appear in *The Journal of Symbolic Logic* 80(3) (2015): 825–865.
- Leigh, G. & C. Nicolai (2013) Axiomatic Truth, Syntax, Metatheoretic Reasoning. *The Review of Symbolic Logic*.
- Leitgeb, H. (2007). What theories of truth should be like (but cannot be). *Blackwell Philosophy Compass* 2/2, 276–290, Blackwell 2007.
- Łeżyk, M. and B. Wcisło (2014). Models of Weak Theories of Truth. Manuscript.
- Lutz, S. What Was the Syntax-Semantics Debate in the Philosophy of Science About? *Philosophy and Phenomenological Research*, forthcoming. DOI:10.1111/phpr.12221.
- McGee, V. (1991). *Truth, Vagueness and Paradox. An Essay on the Logic of Truth*. Hackett, Indianapolis.
- Nicolai, C. (2015). Deflationary Truth and the Ontology of Expressions. *Synthese* Volume 192, Issue 12, pp 4031–4055.
- Nicolai, C. (2016) A Note on typed truth and consistency assertions. *Journal of Philosophical Logic* 45(1): 89–119.
- Nicolai, C. (2016). More on Systems of Truth and Predicative Comprehension. Forthcoming in Boccuni F. and Sereni A. (Eds). *Philosophy of Mathematics: Objectivity, Cognition and Proof*. Boston Studies in the History and Philosophy of Science, Springer.
- Pozsgay, L. J. (1968). Gödel’s second theorem for elementary arithmetic. *Zeitschr. f. math. Logik und Grundlagen d. Math.* 14: 67–80.
- Pudlák, P. (1985). Cuts, Consistency Statements and Interpretations. *The Journal of Symbolic Logic* 50(2): 423–441.
- Rogers, H. (1987). *Theory of Recursive Functions and Effective Computability*. MIT Press, 1987.
- Schichtenberg, H. and S. Wainer. *Proofs and Computations*. ASL Lecture Notes Series, Cambridge University Press.

- Sheard, M. (1994). A guide to truth predicates in the modern era. *Journal of Symbolic Logic* 59, 1032–1054.
- Simpson, S.G. (2009). *Subsystems of Second-Order Arithmetic*. Cambridge University Press.
- Smoryński, C. (1977). The incompleteness theorems. in J. Barwise (ed.), *Handbook of Mathematical Logic*. North Holland.
- Takeuti (1987). *Proof Theory. Second Edition*. North-Holland, Amsterdam.
- Visser, A. (1991). The formalization of interpretability. *Studia Logica* 50, 1, pp: 81–105.
- Visser (1992). An inside view of Exp. *The Journal of Symbolic Logic* 57(1): 131–165.
- Visser, A. (1997). An Overview of Interpretability Logic. In *Advances in Modal Logic '96*.
- Visser, A. (2006). Categories of Theories and Interpretations. In Enayat, A., Kalantari, I. and Moniri, M., (eds). *Logic in Tehran*, Vol. 26. Lecture Notes in Logic, La Jolla, CA: Association for Symbolic Logic, pp. 284–341.
- Visser, A. (2015). The interpretability of inconsistency: Feferman’s theorem and related results. Forthcoming in the *Bulletin of Symbolic Logic*.
- Welch, P. (2001). On Gupta-Belnap Revision Theories, Kripkean Fixed-Points, and the Next Stable Set. *Bulletin of Symbolic Logic* 7(3): 345–360.
- Williamson, T. (201?). Semantic Paradoxes and Abductive Methodology. To appear.
- Woodfield, A. (1991). Conceptions *Mind* 100(4): 547–572.

MUNICH CENTER FOR MATHEMATICAL PHILOSOPHY
 GESCHWISTER-SCHOLL PLATZ 1, MUNICH
 E-mail: Carlo.Nicolai@lrz.uni-muenchen.de